

1 Exploring Genetic Associations of Three Types of Risk Factors with Ischemic Stroke: An
2 Integrated Bioinformatics Study

3

4 Yi Liu^{1#}, PhD; Weili Wang^{1#}, MMedSci; Xin Cui^{1*}, PhD; Yanming Xie^{1*}, Prof.

5

6 ¹Institute of Basic Research in Clinical Medicine, China Academy of Chinese Medical
7 Sciences Beijing, 100700, China;

8

9 [#]Yi Liu and Weili Wang contributed equally to this work.

10

11 ^{*}Yanming Xie and Xin Cui are the joint corresponding authors of this article.

12

13 Correspondence:

14 Yanming Xie No. 16, Nanxiaojie, Dongzhimennei, Dongcheng District, Beijing, China

15 E-mail: ketizu2018@163.com

16 Xin Cui No. 16, Nanxiaojie, Dongzhimennei, Dongcheng District, Beijing, China

17 E-mail: Xinrobertcm@hotmail.com

18

19 Short Title: Associations of exposures and ischemic stroke

20

21 The total number of words: 8817 words

22

23

1 Abstract:

2 Background: Ischemic stroke (IS) is a primary cause of disability and mortality globally.

3 More and more reports suggest a strong association between blood pressure, blood glucose,
4 and blood lipids and their metabolic products with IS.

5 Methods: We extracted the genetic tools of blood pressure, blood glucose, and blood lipids
6 and their metabolites as instrumental variables, which were then paired with GWAS data on
7 IS and a Mendelian randomization (MR) analysis was performed to assess the effect of these
8 exposures on the disease. Following the positive results, colocalization analysis was
9 performed to identify shared genes associated with exposures and IS. We then performed
10 differential expression analysis using the GEO dataset to identify the differentially expressed
11 associated genes (DEAGs) from associated shared genes. Additional analyzes were
12 performed on these DEAGs to obtain their importance scores using four machine learning
13 models. A nomogram was created using genes with high importance scores to predict the
14 level of risk assessment between DEAGs and IS.

15 Results: There is a positive correlation between blood pressure, blood glucose and the risk of
16 IS onset, while blood lipids and their metabolic products are positively or negatively
17 correlated with the risk. There are 64 shared genes of blood pressure, blood lipids and their
18 metabolic products with IS. Thirteen DEAGs were obtained, and among which FURIN,
19 MAN2A2, HDDC3, ALDH2, and TOMM40 were identified as feature genes for creating the
20 nomogram which can quantitatively predict the risk of IS onset with the expression of these
21 feature genes. By cluster analysis, we found that DEAGs expression underlying immune
22 inflammation, angiogenesis and development, lipid metabolism, etc.

1 Conclusion: This study suggests a significant association between blood pressure, blood
2 glucose, and blood lipids and their metabolic products with IS, and predicts that these
3 exposures mainly regulate the occurrence, development, and prognosis of IS through
4 mechanisms such as DNA repair, DNA methylation, mitochondrial repair, apoptosis,
5 autophagy, etc.

6

7 Key words: Mendelian Randomization; Colocalization analysis; Bioinformatics; Ischemic
8 stroke; Blood pressure; Blood glucose; blood lipids and their metabolic products;

9

10 Abbreviations: Area Under the Curve, AUC; Deoxyribonucleic acid, DNA; Differentially
11 expressed associated genes, DEAGs; differentially expressed genes, DEGs; Extreme Gradient
12 Boosting model, XGB; false discovery rate, FDR; Gene Expression Omnibus , GEO; Gene
13 ontology, GO; gene set variation analysis, GSVA; Generalized Linear, GL; genome-wide
14 association study, GWAS; high-density lipoprotein, HDL; Instrumental variables, IVs;
15 Integrated Epidemiology Unit, IEU; intermediate-density lipoprotein, IDL; inverse variance
16 weighted, IVW; Ischemic stroke, IS; Kyoto Encyclopedia of Genes and Genomes, KEGG;
17 low-density lipoprotein, LDL; Mendelian randomization, MR; Natural killer, NK; one-way
18 analysis of variance, ANOVA; posterior probability, PP; Principal Component Analysis, PCA;
19 Random Forest model, RF; receiver operating characteristic curves, ROC; Single nucleotide
20 polymorphisms, SNPs; Single-sample gene set enrichment analysis, ssGSEA; Support Vector
21 Machine model, SVM; very low-density lipoprotein, VLDL

22

23

1 Ischemic stroke (IS) is a disease caused by the blockage of blood vessels in the brain,
2 leading to local cerebral hypoxia and ischemia, resulting in the death of brain cells. This
3 blockage is usually caused by thrombosis or embolism, causing neurological deficits in the
4 brain area with insufficient blood supply.¹⁻³ In 2019, the number of patients who died from
5 stroke worldwide reached 6.55 million, making it the main cause of adult disability and death
6 worldwide.^{4,5} Among newly diagnosed stroke cases, IS accounted for 62.4%.⁶ Therefore,
7 identifying the risk factors for IS is crucial for preventing this disease. Blood pressure, blood
8 glucose, and blood lipids and their metabolic products have received widespread attention as
9 potential factors affecting the risk of IS. A large number of clinical and epidemiological
10 studies have shown that controlling blood pressure, blood glucose, and blood lipids can
11 reduce the risk of IS.⁷⁻⁹ In recent years, with a deeper understanding of the human genome,
12 we have come to realize that the expression of these biomarkers is largely influenced by
13 genetic information.¹⁰ Therefore, exploring high-risk groups from a genetic perspective and
14 identifying shared genes related to blood pressure, blood glucose, and blood lipids and their
15 metabolic products with IS is of great significance for the precise prevention and treatment of
16 IS.

17 Mendelian randomization (MR) is a genetic epidemiological research tool that is widely
18 used to assess the potential causal relationship between exposure and disease.^{11,12} The
19 working mechanism of MR analysis is similar to that of a natural randomized controlled trial.
20 It uses genetic instruments (single nucleotide polymorphisms, SNPs) that are strongly
21 correlated with exposure factors and are not affected by confounding factors under the
22 principle of MR distribution as instrumental variables (IVs).^{13,14} This study used a

1 two-sample MR analysis to explore the potential causal effects of blood pressure, blood
2 glucose, and blood lipids and their metabolic products on IS. The selected datasets are all
3 from the Integrated Epidemiology Unit (IEU) database (<https://gwas.mrcieu.ac.uk/>). This
4 database is public and contains nearly 2.45 trillion genetic associations from 42,334 summary
5 datasets of genome-wide association study (GWAS).

6 Colocalization analysis is a method for assessing whether two correlated traits are driven
7 by the same genetic mechanism. This method often uses the Bayesian model coloc.¹⁵ The
8 coloc model assumes that in each tested region, each trait has at most one association point,
9 and calculates the posterior probability (PP) of all possible association patterns of the two
10 traits through approximate Bayesian factors: 1) H0: no association; 2) H1: only related to trait
11 1; 3) H2: only related to trait 2; 4) H3: related to both trait 1 and trait 2, but through two
12 independent SNPs; 5) H4: related to both trait 1 and trait 2, through a shared SNP. The
13 posterior probabilities of each association pattern are denoted as PP0 to PP4.¹⁵ A higher PP4
14 value (e.g., PP4 > 50%) provides colocalization support, indicating the existence of shared
15 genetic variation between the two traits.¹⁶ This study used colocalization analysis to find
16 shared genes associated with blood pressure, blood glucose, blood lipids and their metabolic
17 products, and IS. These shared genes may provide important clues for revealing the
18 pathogenesis of IS.

19 The Gene Expression Omnibus (GEO) database collects genomic datasets from different
20 species and tissues, which can be used in conjunction with bioinformatics methods to
21 discover evidence of specific gene expression, identify disease prediction factors, and
22 ultimately help us understand the mechanisms by which the genome regulates physiological

1 and pathological states.¹⁷⁻¹⁹ Therefore, based on the associated shared genes of blood pressure,
2 blood glucose, blood lipids and their metabolic products with IS, this study performed
3 differential analysis on normal samples and IS samples in the GEO dataset, and constructed a
4 machine learning model to screen out feature genes among the associated shared genes. On
5 the one hand, it verifies the MR results through real human samples, and on the other hand, it
6 further screens out biomarkers with higher importance, to further explore the mechanisms of
7 blood pressure, blood glucose, blood lipids and their metabolic products on IS.

8

9 1. Study design

10 In this study, aligned with our research objective, we selected GWAS summary datasets
11 for blood pressure, blood glucose, and blood lipids and their metabolic products with IS. We
12 conducted a two-sample MR analysis and chose positive results for further gene
13 colocalization analysis to identify genes associated with the exposure factors and IS.
14 Subsequently, we retrieved two datasets of IS (training and validation sets), encompassing
15 normal and IS groups, from the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>).
16 Differential analysis, correlation analysis, and immune infiltration analysis were carried out
17 on the normal and IS groups of the training set. With the help of machine learning, we
18 derived feature genes related to blood pressure, blood glucose, blood lipids and their
19 metabolic products, and IS. Following this, the results of the machine learning model were
20 validated using the validation set, and two rounds of clustering analysis on the IS group
21 samples were performed to further investigate the functional mechanisms of the
22 exposure-related genes (Figure 1).

1

2 1.1 Data sources

3 In accordance with the goals of this study, we utilized the IEU database, conducting
4 searches with the keywords of "VLDL (very low-density lipoprotein)", "LDL (low-density
5 lipoprotein)", "IDL (intermediate-density lipoprotein)", "HDL (high-density lipoprotein)",
6 "apolipoprotein", "triglyceride", "fatty acid class", "systolic pressure", "diastolic pressure",
7 "blood glucose", "glycosylated hemoglobin", and "IS". This can identify genetic instruments
8 for blood pressure, blood glucose, and blood lipids and their metabolic products as exposure
9 variables. Concurrently, IS-related genetic instruments were selected as outcome variables for
10 an MR analysis, aiming to explore potential causal relationships between exposures and
11 outcomes. By using "ischemic stroke" as the keyword, setting the data type as array-based
12 expression profiles, and specifying the species as homo sapiens, we retrieved samples from
13 the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>) to obtain gene expression and clinical
14 data of IS patients and healthy individuals. Perl code was employed for gene symbol
15 annotation and data adjustment, thus deriving expression levels of genes related to blood
16 pressure, blood glucose, blood lipids and their metabolites, and IS, in both the normal and IS
17 groups.

18

19 1.2 Mendelian Randomization analysis

20 We utilized the "TwoSampleMR" package to conduct a two-sample MR analysis on the
21 relationship between blood pressure, blood glucose, and blood lipids and their metabolic
22 products with IS.²⁰ In the analysis, we set the parameters as follows: 1) The selected genetic

1 instruments need to have a strong correlation with the exposure factors, and the
2 corresponding P-value should be less than 5×10^{-8} ; 2) In conducting Linkage Disequilibrium
3 (LD) clumping, we set the r^2 threshold at 0.01; 3) When performing clumping analysis, we set
4 the window size as 10,000 kilobase pairs.²¹ These stringent parameter settings help us to
5 accurately reveal the causal relationship between exposures and outcomes.

6

7 1.3 Colocalization analysis

8 After completing the two-sample MR analysis, we selected positive results that
9 demonstrated significant associations between the exposure factors and outcome in the MR
10 analysis. We used the "ieugwasr" and "coloc" packages to conduct colocalization analysis on
11 these positive results to determine if the exposure factors and outcome might be associated
12 within the same gene region. We set a posterior probability (PP4) greater than 50% as the
13 standard, indicating that the genes in that region are associated with both exposure and
14 outcomes¹⁶. Finally, we employed the "gassocplot2" package to visualize the colocalization
15 results, clearly demonstrating the associated shared genes of blood pressure, blood glucose,
16 and blood lipids and their metabolic products with IS.

17

18 1.4 Identification of differently expressed associated genes (DEAGs) and analyses of DEAGs

19 Expression levels of genes associated with blood pressure, blood glucose, blood lipids
20 and their metabolic products, and IS were extracted from both the normal and IS groups.
21 Differential expression analysis was performed using the "limma", "pheatmap", "ggpubr" and
22 other R packages, and results were presented as box plots and heatmaps. Genes with a
23 p-value less than 0.05 were defined as differentially expressed associated genes (DEAGs).

1 Using Perl coding, DEAGs were located on chromosomes and their positions were displayed
2 on a circos plot created with the "Rcirco" package. In addition, the "cor" command was used
3 to calculate the correlation coefficient between each two DEAGs and the results were
4 visualized to illustrate the correlation among DEAGs.

5

6 1.5 Analysis of immune cells in IS samples

7 We performed 1000 simulations using the CIBERSORT command in R to obtain a total
8 relative amount of immune cells equal to 1, then visualized the content of immune cells in
9 each sample using a bar plot. Single-sample gene set enrichment analysis (ssGSEA) was
10 performed using the "GSVA" and "GSABase" packages to compare the differences in
11 immune cell contents between the normal group and the IS group, and the results of ssGSEA
12 are presented as box plot. We matched the differentially expressed associated genes (DEAGs)
13 with the ssGSEA scores, performed a correlation test to obtain the correlation coefficient, and
14 then visualized the results as a heat map.

15

16 1.6 Identification of feature genes and construction and validation of nomogram based on ML

17 Expression data from DEAGs was utilized to construct four predictive models: Random
18 Forest model (RF), Support Vector Machine model (SVM), Generalized Linear Model (GL),
19 and Extreme Gradient Boosting model (XGB). The results from these four models were
20 calculated based on the prediction function. Feature genes within DEAGs were then selected
21 through comprehensive analysis of reverse cumulative distribution plots, residual box plots,
22 and receiver operating characteristic (ROC) curves. After selecting the optimal model, line
23 plots were constructed using the expression levels of the feature genes in both the normal and

1 IS groups. Lastly, decision and calibration curves were drawn to assess the accuracy of the
2 line plots. Another dataset that includes both a normal and IS group was obtained from the
3 GEO database, and a machine learning model was constructed using the same R language
4 approach as before. The ROC was plotted to validate the machine learning model constructed
5 in the test dataset.

6

7 1.7 Clustering of DEAGs and analysis between DEAG clusters.

8 The "ConsensusClusterPlus" package in R was used for clustering, employing Euclidean
9 distance type and allowing up to nine clusters. Expression levels between clusters were
10 compared using heatmaps and box plots, and principal component analysis (PCA) was
11 applied to assess differences between clusters. Following this, a ssGSEA was performed on
12 the DEAG clusters, generating bar plots of the amount of individual immune cells in each
13 sample within the different clusters and comparing the differences in immune cell content
14 among different clusters. Gene ontology (GO) and Kyoto Encyclopedia of Genes and
15 Genomes (KEGG) enrichment analyses were performed using gmt files downloaded from the
16 GSEA platform (<http://www.gsea-msigdb.org/>), and gene set variation analysis (GSVA) was
17 conducted in R language to analyze the expression of enrichment items between clusters.
18 Finally, under the filtering conditions of $|\logFC| > 1$ and an adjusted P-value < 0.05 ,
19 differential expression analysis was performed on the gene expression of DEAG clusters.
20 Intersection of DEAG clusters through a Venn diagram yielded differentially expressed genes
21 (DEGs).

22

23 1.8 Clustering of DEGs and analysis between DEG clusters

1 The same clustering method as in section 1.7 was applied to cluster DEGs, and the
2 DEG cluster with the highest precision was selected. Based on the DEGs clustering results,
3 we compared the expression levels of DEGs, the differences in DEAGs expression, and the
4 immune cell content in different clusters. These results were visualized using heatmaps and
5 box plots.

6

7 1.9 Construction of DEAG scores

8 The PCA method was employed to calculate the expression levels of DEAGs for each
9 sample, yielding DEAG scores (The formula is shown below).²² Using R packages such as
10 "limma" and "ggpubr", we performed differential analysis on the scores of DEAGs in both
11 significantly differentially expressed core gene clusters and DEG clusters. We created box
12 plots to illustrate the scores of significantly differentially expressed core genes in samples
13 within the significantly differentially expressed core gene and DEG clusters. In addition, an
14 alluvial diagram was drawn using the package "ggalluvial" to visualize the relationships and
15 overall processes among DEAG clusters, DEG clusters, samples with higher scores, and
16 samples with lower scores of DEAGs.

17

$$\text{DEAG Score} = \sum (\text{PC1}_i + \text{PC2}_i)$$

18

19 1.10 Statistical analysis

20 This study carried out MR analysis and colocalization analysis using R V4.2.1. In the
21 MR analysis, if the number of SNPs for the exposure factor is less than 5 after filtering, we
22 will use the Wald ratio method for analysis; otherwise, we will opt for the inverse variance

1 weighted (IVW) method for analysis. In colocalization analysis, to maximize the acquisition
2 of shared genes between "exposure" and "outcome", we set the chromosomal locus parameter
3 range from 75 kilobase pairs to 500 kilobase pairs, aiming to include a wider area around the
4 SNP to enhance the likelihood of discovering shared genes related to both exposure and
5 outcome. In the bioinformatics part, we used Strawberry Perl 5.32.1.1 for GEO dataset
6 extraction and data annotation, and all other statistical analyses were performed with R
7 V4.2.1. For two independent samples, we applied a t-test, while for two paired samples we
8 used the Wilcoxon paired-rank sum test, and for three or more groups of data, we employed
9 one-way analysis of variance (ANOVA) and the Kruskal-Wallis rank sum test. The Spearman
10 rank correlation test was used for correlation analysis. We set a P value < 0.05 or a false
11 discovery rate (FDR) < 0.05 after correction by the Benjamini-Hochberg method as being
12 statistically significant.

13

14 2 Results

15

16 2.1 Results of collection of GWAS and GEO datasets

17 We retrieved 516 relevant exposure datasets from the ieu database based on the set
18 keywords, and the IS dataset came from 440,328 participants of European ancestry, including
19 34,217 cases and 406,111 controls. These datasets were used in two-sample MR analysis for
20 IS to explore their potential causal relationships. By searching the GEO database for IS
21 datasets, we filtered out two datasets that met the requirements, namely GSE16561 and
22 GSE22255, using the former as the test set and the latter as the validation set. The former is

1 total RNA data from human peripheral blood samples, including 39 IS samples (diagnosed by
2 MRI) and 24 normal samples; the latter is RNA expression data from human peripheral blood
3 monocytes, including 20 IS samples and 20 normal samples.

4

5 2.2 Results of Mendelian randomization analysis

6 Among the 516 exposure datasets, 114 showed a positive association with IS. These
7 positive results further suggested that increased levels of Apolipoprotein A-I and HDL may
8 reduce the risk of IS, while elevated levels of VLDL, LDL, IDL, Apolipoprotein B,
9 triglycerides, fatty acids, systolic blood pressure, diastolic blood pressure, fasting blood sugar,
10 and glycated hemoglobin might increase the risk of IS. (Please see Attachment 1 for detailed
11 analysis results)

12

13 2.3 Results of colocalization analysis

14 Upon completion of the two-sample MR analysis, we further subjected the 114 positive
15 results to colocalization analysis to study the potential shared genes between the exposure
16 factors and IS. The results showed that the PP4 was greater than 50% for 70 of the positive
17 results, including exposure factors such as LDL, IDL, VLDL, Apolipoprotein B, fatty acids,
18 systolic and diastolic blood pressure. This suggests that these exposure factors may have
19 genes shared with IS. The colocalization results were visualized using the "gassocplot2"
20 package, and we finally identified 64 shared genes between these exposure factors and IS.
21 These genes may play a key role in the association between the selected exposure factors and

1 IS. (Please see Attachment 1 for detailed analysis results and Attachment 2 for specific shared
2 genes situations)

3

4 2.4 Results of DEAGs identification and analysis of DEAGs.

5 Through colocalization analysis, we identified 64 associated shared genes between
6 blood pressure, blood sugar, and blood lipids and their metabolites with IS. Differential
7 analysis between the IS group and the normal group in the test dataset revealed that 13
8 associated shared genes exhibited statistically significant differences. Among them, ERCC2,
9 FEN1, HDDC3, TOMM40, TRAPPC6A were highly expressed in the normal group, while
10 ALDH2, FADS2, FES, FURIN, MAN2A2, PVRL2, SYT7, and FAM109A were highly
11 expressed in the IS group, as shown in Figures A and B. The specific positions of these
12 significantly differentially expressed association genes on the chromosome are shown in
13 Figure C. The correlation analysis between each pair of DEAGs in the IS samples showed
14 that there is a certain correlation between DEAGs, which are mainly positive correlations, as
15 shown in Figures D and E.

16

17 2.5 Results of Immune cell infiltration analysis, immune cell difference analysis and
18 correlation analysis of IS samples

19 Through immune cell infiltration analysis, we obtained the types and content of immune cells
20 expressed in each sample, as shown in Figure A. ssGSEA analysis revealed (Figure B) that
21 the expression differences of 11 types of immune cells in the control group and IS group were
22 statistically significant. Among them, B cells naive, T cells CD8, T cells CD4 memory

1 activated, T cells follicular helper, Natural killer (NK) cells activated, and Dendritic cells
2 activated were highly expressed in the control group, while T cells gamma delta,
3 Macrophages M0, Macrophages M2, Mast cells activated, and Neutrophils were highly
4 expressed in the IS group. The correlation analysis between DEAGs and immune cells
5 showed that there is a certain correlation between DEAGs and immune cells, and the levels of
6 positive and negative correlations are comparable, as shown in Figure C.

7

8 2.6 Results of selection of machine learning models, construction of nomogram, and
9 verification

10 The results of building SVM, RF, XGB, and GL machine learning prediction models
11 using DEAGs data were obtained. From the residual boxplots, inverse cumulative distribution
12 plots, and ROC curves (Figure A, C, and D), it can be observed that the GL method had the
13 highest area under the ROC curve, the lowest residual values, and the lowest inverse
14 cumulative values. Therefore, the GL method was considered the most accurate and chosen
15 as the best model for further analysis. The GL model provided importance scores for the
16 selected feature genes, as shown in Figure B, revealing nine feature genes. Among them, the
17 top five genes with the highest importance scores were used to construct the column line
18 graph (including FURIN, MAN2A2, HDDC3, ALDH2, and TOMM40). Then, separate
19 scoring scales were obtained for these five feature genes. The risk rate of the co-located
20 feature genes associated with blood pressure, blood glucose, and blood lipids and their
21 metabolic products in the occurrence of IS was assessed by calculating the sum of the feature
22 gene expression scores (Figure E). The predictive accuracy of the model was evaluated by the

1 distance between the solid line and the dotted line in the calibration curve (Figure F) and the
2 distance between the red and gray lines in the decision curve (Figure G), indicating high
3 accuracy. The GL model was validated using the validation dataset (GSE22255) and the top
4 five feature genes with the highest importance scores were included in the model validation
5 and ROC curve analysis. The results showed an Area Under the Curve (AUC) value of 0.528,
6 with a 95% confidence interval of 0.167-0.889. Therefore, based on the validation of the
7 model, it can be concluded that the model based on the GEO dataset has good accuracy.

8

9 2.7 Results of Clustering of DEAGs, Analysis of DEAGs Expression between Clusters, 10 Immune Cell Analysis, and GO and KEGG Enrichment Analysis

11 Clustering analysis based on the expression of DEAGs revealed two distinct clusters
12 with the highest accuracy, as shown in Figure A. Consequently, the IS samples were divided
13 into two groups: C1 and C2, as depicted in Figures B and C. Subsequently, DEAGs
14 expression analysis was performed between the two DEAG clusters (Figure D and E),
15 indicating that ALDH2, FEN1, FES, and FURIN were significantly upregulated in C1. PCA
16 analysis (Figure F) demonstrated that DEAGs can differentiate between C1 and C2.
17 Moreover, ssGSEA analysis (Figure G) identified three statistically significant immune cell
18 types between C1 and C2. Specifically, NK cells resting and Macrophages M2 were
19 upregulated in C1, while Mast cells resting were upregulated in C2. The respective content of
20 different immune cells in each C1 and C2 sample is illustrated in Figure H.

21 In GSVA analysis (Figures I and J), it can be observed that compared to C1, C2 showed
22 upregulation of several GO biological processes, including Interleukin 6 production, positive

1 regulation of interleukin 6 production, cellular response to osmotic stress, neutrophil
2 degeneration, response to lipoteichoic acid, positive regulation of reactive oxygen species
3 metabolic process, and macropinocytosis. On the other hand, C2 exhibited downregulation of
4 branching involved in blood vessel morphogenesis, cellular response to potassium ion,
5 myotube cell development, regulation of heart morphogenesis, canonical Wnt signaling
6 pathway involved in osteoblast differentiation, gonadotropin secretion, vascular process in
7 the circulatory system, and blood vessel maturation in GO biological processes. In GO
8 molecular functions, C2 showed upregulation of leukotriene C4 synthase activity and
9 downregulation of sphingosine 1-phosphate phosphatase activity. In GO cellular components,
10 C2 exhibited upregulation of early phagosome, specific granule lumen, and
11 proton-transporting two-sector ATPase complex catalytic domain, while downregulation of
12 photoreceptor disc membrane. In KEGG pathways, C2 demonstrated upregulation of other
13 glycan degradation, toll-like receptor signaling pathway, sphingolipid metabolism, pentose
14 phosphate pathway, pathogenic *Escherichia coli* infection, glutathione metabolism, Fc gamma
15 R-mediated phagocytosis, neurotrophin signaling pathway, cell signaling in *Helicobacter*
16 *pylori* infection, progesterone-mediated oocyte maturation, and endocytosis. Conversely, C2
17 showed downregulation of alanine, aspartate, and glutamate metabolism, cytokine-cytokine
18 receptor interaction, taste transduction, ECM receptor interaction, hedgehog signaling
19 pathway, arachidonic acid metabolism, neuroactive ligand-receptor interaction, ascorbate and
20 aldarate metabolism, pentose and glucuronate interconversions, and nitrogen metabolism in
21 KEGG pathways.

22

1 2.8 Results of DEGs screening, DEGs clustering, and analyzes of DEG clusters

2 After filtering for significantly differentially expressed genes (DEGs) between C1 and
3 C2 samples, a total of 4 DEGs were identified (Figure A). These 4 DEGs were further
4 subjected to clustering analysis within the DEAGs cluster, resulting in two distinct clusters,
5 namely CI and CII (Figure B). Analysis of DEAGs expression in these two clusters in IS
6 samples revealed that RNASE2, RNASE3, MMP9, and CAMP exhibited statistically
7 significant differential expression, with higher expression in CI and lower expression in CII
8 (Figure C).

9 The differential expression analysis of DEAGs based on DEGs clustering was
10 performed on the samples, as shown in Figure D. The results revealed that these genes were
11 mainly upregulated in the CI group and downregulated in the CII group. Among them, the
12 significantly differentially expressed genes included ALDH2, FES, MAN2A2 (upregulated in
13 CI), and TRAPPCA (upregulated in CII). Additionally, using ssGSEA based on DEGs
14 clustering (using the same gene set file as above), statistically significant differences in
15 immune cells were observed, including activated B cells (downregulated in CI), activated
16 dendritic cells, immature dendritic cells, myeloid-derived suppressor cells (MDSCs),
17 macrophages, mast cells, monocytes, neutrophils, and plasmacytoid dendritic cells
18 (upregulated in CI), as shown in Figure E.

19

20 2.9 Results of DEAGs scoring, differential analysis of DEAGs score between clusters, and
21 construction of the alluvial plot

22 The differential analysis of DEGs clustering based on the principal component analysis

1 (PCA) scores showed statistically significant differences between the clusters (Figure A).
2 Specifically, C1 had lower scores, while C2 had higher scores. However, there was no
3 statistically significant difference observed in the DEGs clustering (Figure B). The scatter
4 plot in Figure C displayed that the C1 and C2 clusters of DEAGs corresponded mainly to the
5 CI and CII clusters of DEGs, respectively. However, there was no clear correspondence
6 observed between the high and low scores of DEAGs and the DEGs clustering.

7

8 3 Discussion

9

10 3.1. The relationships of blood pressure, blood glucose, and blood lipids and their metabolic
11 products with IS

12 Clinical studies have demonstrated that effective management of various physiological
13 indicators can significantly impact the risk of IS. For example, blood pressure control has
14 been shown to significantly reduce the risk of IS recurrence.²³ Elevated fasting blood glucose
15 levels are significantly associated with an increased risk of IS.²⁴ Furthermore, increased
16 levels of glycated hemoglobin, regardless of diabetes status, are associated with an increased
17 risk of IS.²⁵ Lipid management also plays a significant role in IS risk. There are significant
18 differences in VLDL levels between the control and IS groups.²⁶ However, the role and
19 mechanisms of VLDL level changes in IS pathogenesis require further investigation.
20 Medications that lower LDL levels can significantly reduce the risk of IS.²⁷ Increased levels
21 of IDL and apolipoprotein B are associated with an increased risk of IS.^{28,29} Controlling
22 triglyceride levels can lower the risk of IS.³⁰ Both polyunsaturated fatty acids and saturated

1 fatty acids are also associated with IS.^{31,32} However, elevated levels of HDL and
2 apolipoprotein A-I are associated with a lower risk of IS.^{33,34} Taking all the above clinical
3 studies into consideration, it can be concluded that blood pressure, blood glucose, and blood
4 lipids and their metabolic products play critical roles in influencing the risk of IS.

5 The results of MR analysis indicate a potential causal relationship between blood
6 pressure (diastolic and systolic), blood glucose (fasting blood glucose and glycated
7 hemoglobin), blood lipids and their various metabolic products (VLDL, LDL, IDL, HDL,
8 apolipoprotein A-I, apolipoprotein B, triglycerides, Eicosapentaenoate, Docosapentaenoate,
9 Stearidonate, Docosahexaenoic acid) with IS. Increasing levels of high-density lipoprotein
10 (HDL) and apolipoprotein A-I are associated with a decreased risk of IS. However, elevated
11 blood pressure (including systolic and diastolic), blood glucose (fasting blood glucose and
12 glycated hemoglobin), various low-density lipoproteins such as VLDL, LDL, IDL, as well as
13 apolipoprotein B, triglycerides, Eicosapentaenoate, Docosapentaenoate, Stearidonate,
14 Docosahexaenoic acid are associated with an increased risk of IS.

15

16 3.2 Discussion of mechanisms associated shared genes

17 Through MR analysis, we found that blood pressure, blood glucose, blood lipids and
18 their metabolic products have a significant impact on the risk of IS. To further explore the
19 specific mechanisms underlying these effects, we conducted a colocalization analysis. The
20 results showed that there are 64 associated shared genes between LDL, IDL, VLDL,
21 apolipoprotein B, fatty acids, systolic blood pressure, and diastolic blood pressure with IS.
22 Although a positive association between blood glucose and IS risk has been established, we

1 did not identify any associated shared genes due to the PP4 being less than 50%.

2 The differential analysis of the 64 associated shared genes between normal and IS
3 samples, using the GEO dataset, revealed 13 genes with significant differential expression.
4 Among them, ERCC2, a transcription factor (TF) involved in nucleotide excision repair of
5 damaged DNA,³⁵ has been clinically associated with increased stroke risk,^{36,37} and animal
6 studies have shown its neuroprotective role in preventing ischemia-reperfusion injury.³⁸
7 FEN1, another TF encoding a protein involved in rDNA and mitochondrial DNA repair,³⁹ has
8 been implicated in the pathogenesis of IS and is considered a therapeutic target for IS
9 treatment.^{40,41} HDDC3, a protein encoded by HDDC3 gene, participates in the starvation
10 response, and starvation-induced autophagy is known to protect neurons and regulate IS.^{42,43}
11 TOMM40, encoding the translocase of the outer mitochondrial membrane 40 homolog, has
12 been identified as a potential susceptibility locus for IS based on a study conducted in
13 Japan.⁴⁴ ALDH2, encoding aldehyde dehydrogenase, has shown neuroprotective effects by
14 clearing 4-hydroxy-2-nonenal and reducing mitochondrial-associated cell apoptosis through
15 JNK-mediated cystathionine- β -synthase-3 activation, making it a potential target for IS
16 intervention. Case-control studies have also suggested an association between ALDH2
17 polymorphisms and IS risk in the Han Chinese population.⁴⁵⁻⁴⁷ FADS2, encoding fatty acid
18 desaturase 2, has been associated with IS risk and lipid levels in case-control studies,
19 although the specific mechanisms are still under investigation.⁴⁸ The product of the FES gene
20 exhibits tyrosine-specific protein kinase activity. Inhibitors of this class of enzymes have
21 shown potential therapeutic value for various diseases, including IS, but the involvement of
22 FES in IS is yet to be confirmed.⁴⁹ FURIN, encoding proprotein convertase subtilisin/kexin

1 type 1, has been found to be upregulated within 24 hours after cerebral ischemia in animal
2 experiments,⁵⁰ and case-control studies have revealed lower DNA methylation levels in the
3 IS region, suggesting its association with IS.⁵¹ SYT7, encoding synaptotagmin 7, has been
4 identified to have significant DNA methylation correlation with IS in a large-scale
5 sequencing study.⁵² The roles of FAM109A, MAN2A2, PVRL2, and others in IS require
6 further investigation.

7 These genes primarily function through DNA repair, DNA methylation, mitochondrial
8 repair, cell apoptosis, and autophagy to regulate lipid levels, blood pressure, and neuronal
9 protection, thereby affecting the occurrence, development, and prognosis of IS. The
10 associated shared genes are predominantly upregulated in the IS group, and they exhibit
11 mostly positive regulatory relationships, indicating potential synergistic effects among these
12 genes in IS. Differential expression analysis of immune cells between the IS and normal
13 groups reveals significant expression differences in 50% of the immune cells, suggesting the
14 involvement of immune cells in the regulation of IS by these associated shared genes.
15 Moreover, the comparable levels of positive and negative regulation of immune cells by
16 DEAGs suggest that DEAGs may have bidirectional regulatory effects on the immune
17 system's role in IS.

18 The feature genes obtained from the GL model constructed through machine learning
19 (FURIN, MAN2A2, HDDC3, ALDH2, TOMM40) have significant importance in regulating
20 the associated shared genes and the risk of developing IS. The column plot provides a
21 quantitative prediction of the aforementioned importance and risk of disease. Furthermore,

1 the construction of the GL model using the validation dataset indicates that the column plot
2 has high accuracy.

3 The IS samples were clustered into C1 and C2 groups based on the expression of
4 DEAGs, with C1 showing predominantly high expression of DEAGs and C2 showing low
5 expression. The two clusters were significantly distinguishable according to the PCA results.
6 Furthermore, immune cell infiltration analysis of DEAGs clustering revealed significant
7 differences in immune cell expression between C1 and C2 groups of IS patients. The GSVA
8 results indicated significant expression differences between C1 and C2 groups in various
9 biological processes and pathways, including immune inflammation, vascular development
10 and formation, reactive oxygen species metabolism, carbohydrate metabolism, lipid
11 metabolism, fatty acid metabolism, amino acid metabolism, and neural regulation.

12 Differential analysis between the C1 and C2 groups identified 4 DEGs, which have been
13 extensively studied and shown to be associated with IS incidence, IS prognosis, and IS
14 complications.⁵³⁻⁵⁵ Subsequently, IS samples were further clustered into CI and CII groups
15 based on the expression of DEGs. Differential analysis revealed that DEGs were highly
16 expressed in the CI group and lowly expressed in the CII group, while DEAGs were mainly
17 upregulated in the CI group and downregulated in the CII group. Immune cell infiltration
18 analysis of DEGs clustering demonstrated significant differences in immune cell expression
19 between the CI and CII groups of IS patients, with upregulation observed in the CI group and
20 downregulation observed in the CII group. PCA-based DEAG scoring and subsequent
21 differential analysis revealed significant differences in DEAG scoring between the C1 and C2
22 groups, while no significant differences were observed between the CI and CII groups.

1 Furthermore, the construction of Sankey diagrams provided a visual representation of the
2 correspondence between the two clustering results and the high/low DEAGs\ scoring groups,
3 with C1 and C2 corresponding to CI, CII, and the high/low DEAG groups, while the
4 correspondence between CI and CII was less clear.

6 3.3 Limitations of this study

7 Although MR analysis and colocalization analysis provide valuable insights into the risk
8 factors and underlying mechanisms of IS from the perspective of exposure factors and genes,
9 these methods also have their limitations. The pathogenesis of IS is extremely complex and
10 involves intricate interactions among multiple environmental factors, physiological indicators,
11 and genetic factors. Therefore, in order to comprehensively understand and elucidate the risk
12 factors and mechanisms of IS, a holistic and multifactorial research approach is needed.
13 Additionally, the GWAS data included in our study primarily originate from sample
14 populations of European ancestry. Hence, extrapolating these results directly to non-European
15 populations may have limitations. Similarly, bioinformatics methods also have certain
16 limitations. Although they have provided multi-faceted and multi-dimensional analyses of
17 human blood samples, elucidating specific mechanisms of interaction between blood pressure,
18 blood glucose, blood lipids and their metabolic products in relation to IS, further research and
19 clinical applications require careful consideration. Moreover, the blood samples used in this
20 study were only collected from the United States and Portugal, and whether the results can be
21 extrapolated to other countries, ethnicities, and populations remains to be investigated and
22 confirmed. Given the relatively small sample size, there is a greater risk of bias in the
23 conclusions. Therefore, large-scale studies or animal experiments can be conducted to further

1 validate and deepen the conclusions.

2

3 4 Conclusion and prospects

4 This study first integrated MR analysis, colocalization analysis, and bioinformatics
5 analysis to conduct an in-depth study of the association of blood pressure, blood glucose, and
6 blood lipids and their metabolic products with IS. Firstly, the results of MR analysis showed
7 that there is a positive correlation between blood pressure, blood glucose and the risk of IS
8 onset, and blood lipids and their metabolic products have either a positive or negative
9 correlation with the risk of IS onset. Furthermore, through colocalization analysis, we
10 identified 64 associated shared genes with blood pressure, blood lipids and their metabolic
11 products in relation to IS. Additionally, using the GEO dataset, we identified 13 DEAGs.
12 These genes primarily function through deoxyribonucleic acid (DNA) repair, DNA
13 methylation, mitochondrial repair, cell apoptosis, and autophagy, regulating lipid metabolism,
14 blood pressure, and neuroprotection, thereby influencing the occurrence, progression, and
15 prognosis of IS. Subsequently, through the construction of machine learning models, we
16 selected feature genes. Cluster analysis revealed that the expression of DEAGs in IS may also
17 involve mechanisms such as immune inflammation, vascular development, lipid metabolism,
18 and fatty acid metabolism.

19 The findings of this study will contribute to a deeper understanding of the
20 pathophysiological mechanisms underlying IS, optimize treatment strategies, and drive the
21 development of new drugs. The methodological approach used in this study, which combines
22 genetic epidemiology with bioinformatics, offers a unique perspective for elucidating causal

1 relationships between exposure factors and diseases, as well as unraveling genetic
2 mechanisms. This research strategy, along with the identification of the associated shared
3 genes, provides a new theoretical framework and practical approach for the prevention,
4 diagnosis, and treatment of IS.

5

6 Article Information

7

8 Acknowledgments

9 The authors would like to thank the institutions for providing GWAS summary statistics and
10 human gene expression datasets.

11

12 Source of Funding

13 National Administration of Traditional Chinese Medicine 2021 Qihuang Scholar Support
14 Project (National TCM Education Letter 2022 No. 6) ; Yanming Xie's National Renowned
15 Traditional Chinese Medicine Expert Inheritance Studio Construction Project (National TCM
16 Education Letter 2022 No. 75).

17

18 Disclosures

19 The authors declare that the research was conducted in the absence of any commercial or
20 financial relationships that could be construed as a potential conflict of interest.

21

22 References:

- 1 1. Hankey GJ. Stroke. *Lancet*. 2017;389:641-654. doi: 10.1016/S0140-6736(16)30962-X.
- 2 2. Lee RHC, Lee MHH, Wu CYC, Couto E Silva A, Possoit HE, Hsieh TH, Minagar A, Lin
3 HW. Cerebral ischemia and neuroregeneration. *Neural Regen Res*. 2018;13:373-385. doi:
4 10.4103/1673-5374.228711.
- 5 3. Powers WJ, Rabinstein AA, Ackerson T, Adeoye OM, Bambakidis NC, Becker K, Biller J,
6 Brown M, Demaerschalk BM, Hoh B, et al. Guidelines for the Early Management of
7 Patients With Acute Ischemic Stroke: 2019 Update to the 2018 Guidelines for the Early
8 Management of Acute Ischemic Stroke: A Guideline for Healthcare Professionals From the
9 American Heart Association/American Stroke Association. *Stroke*. 2019;50:e344-e418. doi:
10 10.1161/STR.0000000000000211.
- 11 4. Romano JG, Rundek T. Expanding Treatment for Acute Ischemic Stroke beyond
12 Revascularization. *N Engl J Med*. 2023;388:2095-2096. doi: 10.1056/NEJMe2303184.
- 13 5. Roth GA, Mensah GA, Johnson CO, Addolorato G, Ammirati E, Baddour LM, Barengo
14 NC, Beaton AZ, Benjamin EJ, Benziger CP, et al. Global Burden of Cardiovascular
15 Diseases and Risk Factors, 1990-2019: Update From the GBD 2019 Study. *J Am Coll*
16 *Cardiol*. 2020;76:2982-3021. doi: 10.1016/j.jacc.2020.11.010.
- 17 6. GBD 2019 Stroke Collaborators. Global, regional, and national burden of stroke and its
18 risk factors, 1990-2019: a systematic analysis for the Global Burden of Disease Study 2019.
19 *Lancet Neurol*. 2021;20:795-820. doi: 10.1016/S1474-4422(21)00252-0.
- 20 7. Hansson L, Zanchetti A, Carruthers SG, Dahlöf B, Elmfeldt D, Julius S, Ménard J, Rahn
21 KH, Wedel H, Westerling S. Effects of intensive blood-pressure lowering and low-dose
22 aspirin in patients with hypertension: principal results of the Hypertension Optimal

- 1 Treatment (HOT) randomised trial. HOT Study Group. *Lancet*. 1998;351:1755-1762. doi:
2 10.1016/s0140-6736(98)04311-6.
- 3 8. Duckworth W, Abaira C, Moritz T, Reda D, Emanuele N, Reaven PD, Zieve FJ, Marks J,
4 Davis SN, Hayward R, et al. Glucose control and vascular complications in veterans with
5 type 2 diabetes. *N Engl J Med*. 2009;360:129-139. doi: 10.1056/NEJMoa0808431.
- 6 9. Amarenco P, Labreuche J. Lipid management in the prevention of stroke: review and
7 updated meta-analysis of statins for stroke prevention. *Lancet Neurol*. 2009;8:453-463. doi:
8 10.1016/S1474-4422(09)70058-4.
- 9 10. Kim YS, Leventhal BL. Genetic epidemiology and insights into interactive genetic and
10 environmental effects in autism spectrum disorders. *Biol Psychiatry*. 2015;77:66-74. doi:
11 10.1016/j.biopsych.2014.11.001.
- 12 11. Smith GD, Ebrahim S. 'Mendelian randomization': can genetic epidemiology contribute to
13 understanding environmental determinants of disease? *Int J Epidemiol*. 2003;32:1-22. doi:
14 10.1093/ije/dyg070.
- 15 12. Grover S, Del Greco M F, Stein CM, Ziegler A. Mendelian Randomization. *Methods Mol*
16 *Biol*. 2017;1666:581-628. doi: 10.1007/978-1-4939-7274-6_29.
- 17 13. Verduijn M, Siegerink B, Jager KJ, Zoccali C, Dekker FW. Mendelian randomization: use
18 of genetics to enable causal inference in observational studies. *Nephrol Dial Transplant*.
19 2010;25:1394-1398. doi: 10.1093/ndt/gfq098.
- 20 14. Davies NM, Holmes MV, Davey Smith G. Reading Mendelian randomisation studies: a
21 guide, glossary, and checklist for clinicians. *BMJ*. 2018;362:k601. doi: 10.1136/bmj.k601.
- 22 15. Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, Plagnol

- 1 V. Bayesian test for colocalisation between pairs of genetic association studies using
2 summary statistics. *PLoS Genet.* 2014;10:e1004383. doi: 10.1371/journal.pgen.1004383.
- 3 16. Yarmolinsky J, Bouras E, Constantinescu A, Burrows K, Bull CJ, Vincent EE, Martin RM,
4 Dimopoulou O, Lewis SJ, Moreno V, et al. Genetically proxied glucose-lowering drug
5 target perturbation and risk of cancer: a Mendelian randomisation analysis. *Diabetologia.*
6 2023. doi: 10.1007/s00125-023-05925-4.
- 7 17. Barrett T, Troup DB, Wilhite SE, Ledoux P, Rudnev D, Evangelista C, Kim IF, Soboleva
8 A, Tomashevsky M, Marshall KA, et al. NCBI GEO: archive for high-throughput
9 functional genomic data. *Nucleic Acids Res.* 2009;37(Database issue):D885-890. doi:
10 10.1093/nar/gkn764.
- 11 18. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA,
12 Phillippy KH, Sherman PM, Holko M, et al. NCBI GEO: archive for functional genomics
13 data sets--update. *Nucleic Acids Res.* 2013;41(Database issue):D991-995. doi:
14 10.1093/nar/gks1193.
- 15 19. Xu B, Wang L, Zhan H, Zhao L, Wang Y, Shen M, Xu K, Li L, Luo X, Zhou S, et al.
16 Investigation of the Mechanism of Complement System in Diabetic Nephropathy via
17 Bioinformatics Analysis. *J Diabetes Res.* 2021;2021:5546199. doi: 10.1155/2021/5546199.
- 18 20. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, Laurin C, Burgess S,
19 Bowden J, Langdon R, et al. The MR-Base platform supports systematic causal inference
20 across the human phenome. *Elife.* 2018;7:e34408. doi: 10.7554/eLife.34408.
- 21 21. Xiang M, Wang Y, Gao Z, Wang J, Chen Q, Sun Z, Liang J, Xu J. Exploring causal
22 correlations between inflammatory cytokines and systemic lupus erythematosus: A

- 1 Mendelian randomization. *Front Immunol.* 202;13:985729. doi:
2 10.3389/fimmu.2022.985729.
- 3 22. Zhang B, Wu Q, Li B, Wang D, Wang L, Zhou YL. m6A regulator-mediated methylation
4 modification patterns and tumor microenvironment infiltration characterization in gastric
5 cancer. *Mol Cancer.* 2020;19:53. doi: 10.1186/s12943-020-01170-0.
- 6 23. Biffi A, Anderson CD, Battey TW, Ayres AM, Greenberg SM, Viswanathan A, Rosand J.
7 Association Between Blood Pressure Control and Risk of Recurrent Intracerebral
8 Hemorrhage. *JAMA.* 2015;314:904-912. doi: 10.1001/jama.2015.10082.
- 9 24. Boden-Albala B, Cammack S, Chong J, Wang C, Wright C, Rundek T, Elkind MS, Paik
10 MC, Sacco RL. Diabetes, fasting glucose levels, and risk of ischemic stroke and vascular
11 events: findings from the Northern Manhattan Study (NOMAS). *Diabetes Care.*
12 2008;31:1132-1137. doi: 10.2337/dc07-0797.
- 13 25. Nomani AZ, Nabi S, Ahmed S, Iqbal M, Rajput HM, Rao S. High HbA1c is associated
14 with higher risk of ischaemic stroke in Pakistani population without diabetes. *Stroke Vasc*
15 *Neurol.* 2016;1:133-139. doi: 10.1136/svn-2016-000018.
- 16 26. Bansal BC, Sood AK, Bansal CB. Familial hyperlipidemia in stroke in the young. *Stroke.*
17 1986;17:1142-1145. doi: 10.1161/01.str.17.6.1142.
- 18 27. Schwartz GG, Steg PG, Szarek M, Bhatt DL, Bittner VA, Diaz R, Edelberg JM,
19 Goodman SG, Hanotin C, Harrington RA, et al. Alirocumab and Cardiovascular Outcomes
20 after Acute Coronary Syndrome. *N Engl J Med.* 2018;379:2097-2107. doi:
21 10.1056/NEJMoa1801174.

- 1 28. Berger JS, McGinn AP, Howard BV, Kuller L, Manson JE, Otvos J, Curb JD, Eaton CB,
2 Kaplan RC, Lynch JK, et al. Lipid and lipoprotein biomarkers and the risk of ischemic
3 stroke in postmenopausal women. *Stroke*. 2012;43:958-66. doi:
4 10.1161/STROKEAHA.111.641324.
- 5 29. Johannesen CDL, Mortensen MB, Langsted A, Nordestgaard BG. ApoB and Non-HDL
6 Cholesterol Versus LDL Cholesterol for Ischemic Stroke Risk. *Ann Neurol*.
7 2022;92:379-389. doi: 10.1002/ana.26425.
- 8 30. Das Pradhan A, Glynn RJ, Fruchart JC, MacFadyen JG, Zaharris ES, Everett BM,
9 Campbell SE, Oshima R, Amarenco P, Blom DJ, et al. Triglyceride Lowering with
10 Pemafibrate to Reduce Cardiovascular Risk. *N Engl J Med*. 2022 ;387:1923-1934. doi:
11 10.1056/NEJMoa2210645.
- 12 31. Saber H, Yakoob MY, Shi P, Longstreth WT Jr, Lemaitre RN, Siscovick D, Rexrode KM,
13 Willett WC, Mozaffarian D. Omega-3 Fatty Acids and Incident Ischemic Stroke and Its
14 Atherothrombotic and Cardioembolic Subtypes in 3 US Cohorts. *Stroke*.
15 2017;48:2678-2685. doi: 10.1161/STROKEAHA.117.018235.
- 16 32. Venø SK, Schmidt EB, Jakobsen MU, Lundbye-Christensen S, Bach FW, Overvad K.
17 Substitution of Linoleic Acid for Other Macronutrients and the Risk of Ischemic Stroke.
18 *Stroke*. 2017;48:3190-3195. doi: 10.1161/STROKEAHA.117.017935.
- 19 33. Curb JD, Abbott RD, Rodriguez BL, Masaki KH, Chen R, Popper JS, Petrovitch H, Ross
20 GW, Schatz IJ, Belleau GC, et al. High density lipoprotein cholesterol and the risk of
21 stroke in elderly men: the Honolulu heart program. *Am J Epidemiol*. 2004;160:150-157.
22 doi: 10.1093/aje/kwh177.

- 1 34. O'Donnell MJ, McQueen M, Sniderman A, Pare G, Wang X, Hankey GJ, Rangarajan S,
2 Chin SL, Rao-Melacini P, Ferguson J, et al. Association of Lipids, Lipoproteins, and
3 Apolipoproteins with Stroke Subtypes in an International Case Control Study
4 (INTERSTROKE). *J Stroke*. 2022;24:224-235. doi: 10.5853/jos.2021.02152.
- 5 35. National Center for Biotechnology Information (2023). PubChem Gene Summary for
6 Gene 2068, ERCC2 - ERCC excision repair 2, TFIIH core complex helicase subunit
7 (human). Retrieved June 14, 2023
8 from <https://pubchem.ncbi.nlm.nih.gov/gene/ERCC2/human>.
- 9 36. Shyu HY, Shieh JC, Ji-Ho L, Wang HW, Cheng CW. Polymorphisms of DNA repair
10 pathway genes and cigarette smoking in relation to susceptibility to large artery
11 atherosclerotic stroke among ethnic Chinese in Taiwan. *J Atheroscler Thromb*.
12 2012;19:316-325. doi: 10.5551/jat.10967.
- 13 37. LD, Dawsey SM, Dong ZW, Taylor PR, Mark SD. A prospective study of polymorphisms
14 of DNA repair genes XRCC1, XPD23 and APE/ref-1 and risk of stroke in Linxian, China.
15 *J Epidemiol Community Health*. 2007;61:737-741. doi: 10.1136/jech.2006.048934.
- 16 38. Zhang J, Guo F, Zhou R, Xiang C, Zhang Y, Gao J, Cao G, Yang H. Proteomics and
17 transcriptome reveal the key transcription factors mediating the protection of Panax
18 notoginseng saponins (PNS) against cerebral ischemia/reperfusion injury. *Phytomedicine*.
19 2021;92:153613. doi: 10.1016/j.phymed.2021.153613.
- 20 39. National Center for Biotechnology Information (2023). PubChem Gene Summary for
21 Gene 2237, FEN1 - flap structure-specific endonuclease 1 (human). Retrieved June 15,
22 2023 from <https://pubchem.ncbi.nlm.nih.gov/gene/FEN1/human>.

- 1 40. He Z, Ning N, Zhou Q, Khoshnam SE, Farzaneh M. Mitochondria as a therapeutic target
2 for ischemic stroke. *Free Radic Biol Med.* 2020;146:45-58. doi:
3 10.1016/j.freeradbiomed.2019.11.005.
- 4 41. Andrabi SS, Parvez S, Tabassum H. Ischemic stroke and mitochondria: mechanisms and
5 targets. *Protoplasma.* 2020;257:335-343. doi: 10.1007/s00709-019-01439-2.
- 6 42. Sun D, Lee G, Lee JH, Kim HY, Rhee HW, Park SY, Kim KJ, Kim Y, Kim BY, Hong JI,
7 et al. A metazoan ortholog of SpoT hydrolyzes ppGpp and functions in starvation
8 responses. *Nat Struct Mol Biol.* 2010;17:1188-94. doi: 10.1038/nsmb.1906.
- 9 43. He S, Wang C, Dong H, Xia F, Zhou H, Jiang X, Pei C, Ren H, Li H, Li R, et al.
10 Immune-related GTPase M (IRGM1) regulates neuronal autophagy in a mouse model of
11 stroke. *Autophagy.* 2012;8:1621-1627. doi: 10.4161/auto.21561.
- 12 44. Yamase Y, Horibe H, Ueyama C, Fujimaki T, Oguri M, Kato K, Arai M, Watanabe S,
13 Yamada Y. Association of *TOMM40* and *SLC22A4* polymorphisms with ischemic stroke.
14 *Biomed Rep.* 2015;3:491-498. doi: 10.3892/br.2015.457.
- 15 45. Guo JM, Liu AJ, Zang P, Dong WZ, Ying L, Wang W, Xu P, Song XR, Cai J, Zhang SQ,
16 et al. ALDH2 protects against stroke by clearing 4-HNE. *Cell Res.* 2013;23:915-930. doi:
17 10.1038/cr.2013.69.
- 18 46. Xia P, Zhang F, Yuan Y, Chen C, Huang Y, Li L, Wang E, Guo Q, Ye Z. ALDH 2
19 conferred neuroprotection on cerebral ischemic injury by alleviating mitochondria-related
20 apoptosis through JNK/caspase-3 signing pathway. *Int J Biol Sci.* 2020;16:1303-1323. doi:
21 10.7150/ijbs.38962.
- 22 47. Sun S, He J, Zhang Y, Xiao R, Yan M, Ren Y, Zhu Y, Jin T, Xia Y. Genetic

- 1 polymorphisms in the ALDH2 gene and the risk of ischemic stroke in a Chinese han
2 population. *Oncotarget*. 2017;8:101936-101943. doi: 10.18632/oncotarget.21803.
- 3 48. Yang Q, Yin RX, Cao XL, Wu DF, Chen WX, Zhou YJ. Association of two
4 polymorphisms in the FADS1/FADS2 gene cluster and the risk of coronary artery disease
5 and ischemic stroke. *Int J Clin Exp Pathol*. 2015;8:7318-7331.
- 6 49. Gaęało I, Rusiecka I, Kocić I. Tyrosine Kinase Inhibitor as a new Therapy for Ischemic
7 Stroke and other Neurologic Diseases: is there any Hope for a Better Outcome? *Curr*
8 *Neuropharmacol*. 2015;13:836-844. doi: 10.2174/1570159x13666150518235504.
- 9 50. Yokota N, Uchijima M, Nishizawa S, Namba H, Koide Y. Identification of differentially
10 expressed genes in rat hippocampus after transient global cerebral ischemia using
11 subtractive cDNA cloning based on polymerase chain reaction. *Stroke*. 2001;32:168-174.
12 doi: 10.1161/01.str.32.1.168.
- 13 51. Peng H, Fan Y, Li J, Zheng X, Zhong C, Zhu Z, He Y, Zhang M, Zhang Y. DNA
14 Methylation of the Natriuretic Peptide System Genes and Ischemic Stroke: Gene-Based
15 and Gene Set Analyses. *Neurol Genet*. 2022;8:e679. doi:
16 10.1212/NXG.0000000000000679.
- 17 52. Zhang H, Mo X, Wang A, Peng H, Guo D, Zhong C, Zhu Z, Xu T, Zhang Y. Association
18 of DNA Methylation in Blood Pressure-Related Genes With Ischemic Stroke Risk and
19 Prognosis. *Front Cardiovasc Med*. 2022;9:796245. doi: 10.3389/fcvm.2022.796245.
- 20 53. Sundström J, Söderholm M, Borné Y, Nilsson J, Persson M, Östling G, Melander O,
21 Orho-Melander M, Engström G. Eosinophil Cationic Protein, Carotid Plaque, and
22 Incidence of Stroke. *Stroke*. 2017;48:2686-2692. doi: 10.1161/STROKEAHA.117.018450.

1 54. Zhong C, Yang J, Xu T, Xu T, Peng Y, Wang A, Wang J, Peng H, Li Q, Ju Z, et al. Serum
2 matrix metalloproteinase-9 levels and prognosis of acute ischemic stroke. *Neurology*.
3 2017;89:805-812. doi: 10.1212/WNL.0000000000004257.

4 55. Ponsaerts L, Alders L, Schepers M, de Oliveira RMW, Prickaerts J, Vanmierlo T,
5 Bronckaers A. Neuroinflammation in Ischemic Stroke: Inhibition of cAMP-Specific
6 Phosphodiesterases (PDEs) to the Rescue. *Biomedicines*. 2021;9:703. doi:
7 10.3390/biomedicines9070703.

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

1 TABLE 1| Gene name conversion table

Gene description	Gene symbol
ERCC excision repair 2, TFIIH core complex helicase subunit	ERCC2
flap structure-specific endonuclease	FEN1
HD domain containing 3	HDDC3
translocase of outer mitochondrial membrane 40	TOMM40
trafficking protein particle complex subunit 6A	TRAPPC6A
aldehyde dehydrogenase 2 family member	ALDH2
fatty acid desaturase 2	FADS2
FES proto-oncogene, tyrosine kinase	FES
furin, paired basic amino acid cleaving enzyme	FURIN
mannosidase alpha class 2A member 2	MAN2A2
nectin cell adhesion molecule 2	PVRL2
synaptotagmin 7	SYT7
PH domain containing endocytic trafficking adaptor 1	FAM109A

2

3

4

5

6

7

8

9

10

11

12

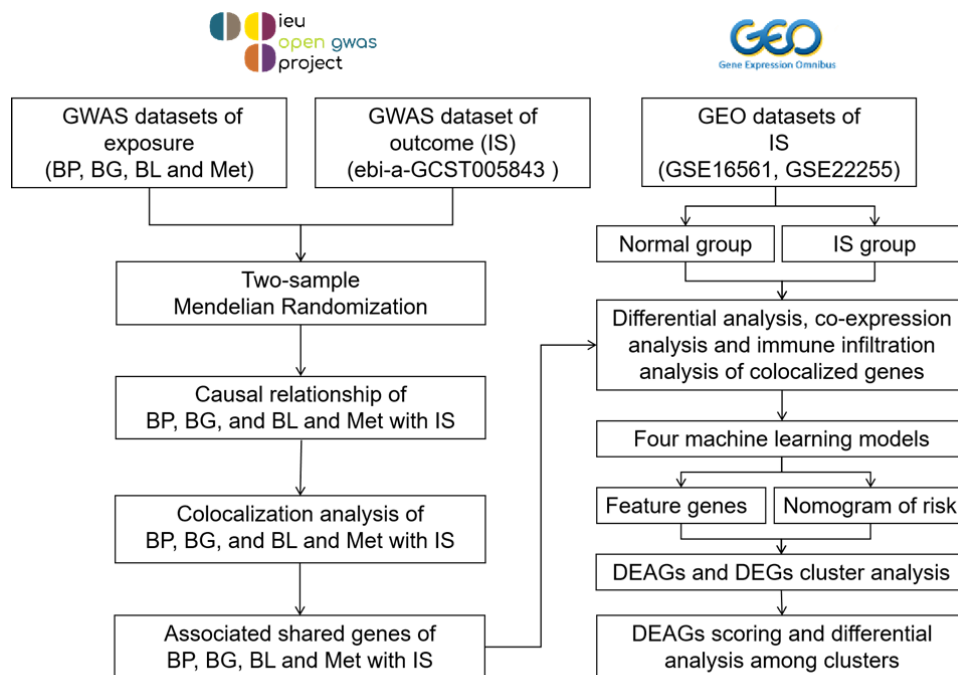
13

14

15

16

1 FIGURE 1 | Flow diagram of this study design



2

3

4

5

6

7

8

9

10

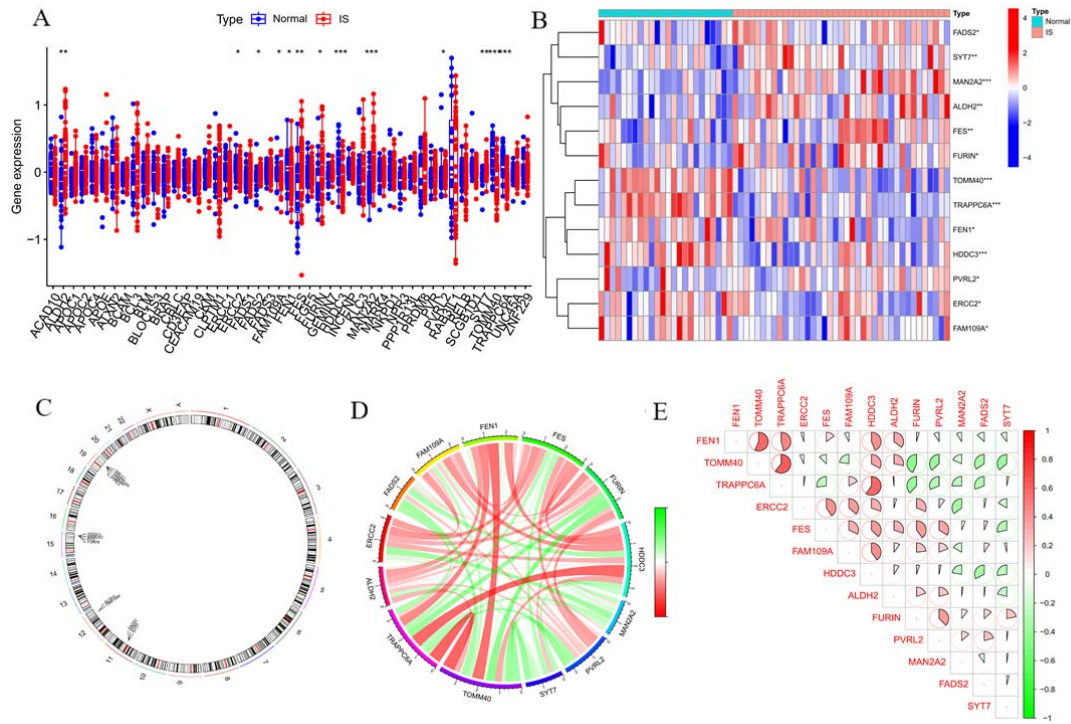
11

12

13

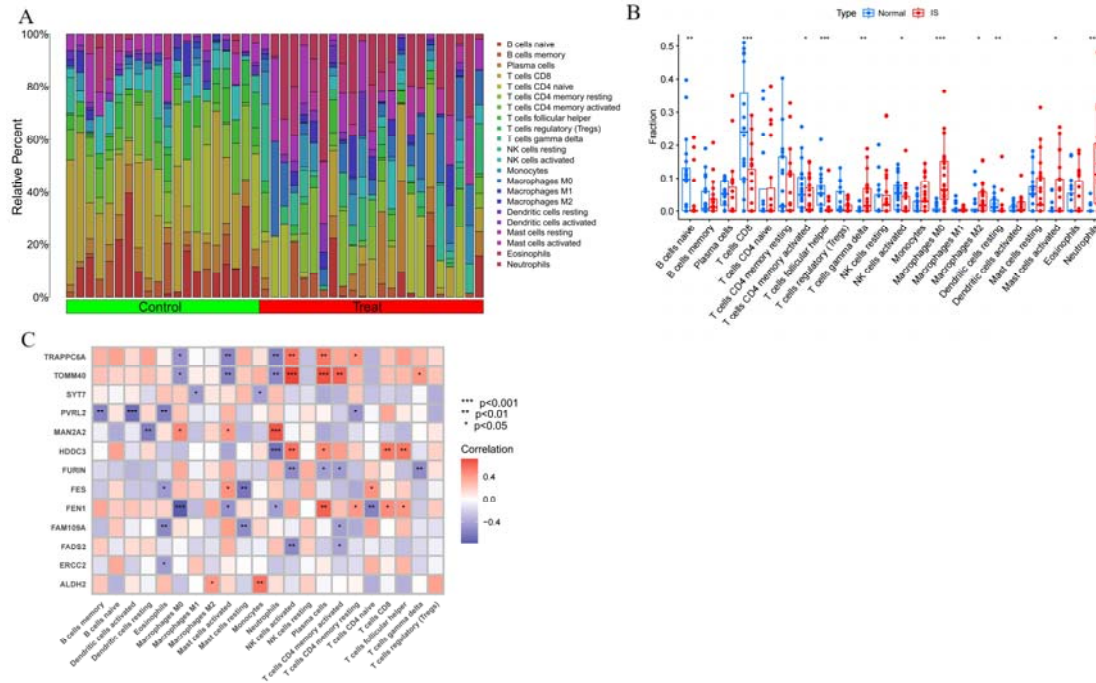
14

- 1 FIGURE 2 | (A) Box plot of expression difference analysis of associated shared genes
- 2 between normal samples and IS samples; (B) Heat map of DEAGs expression in normal and
- 3 IS samples; (C) Circle plot of chromosome location of DEAGs; (D) DEAGs correlation
- 4 network; (E) Correlation analysis between the two DEAGs



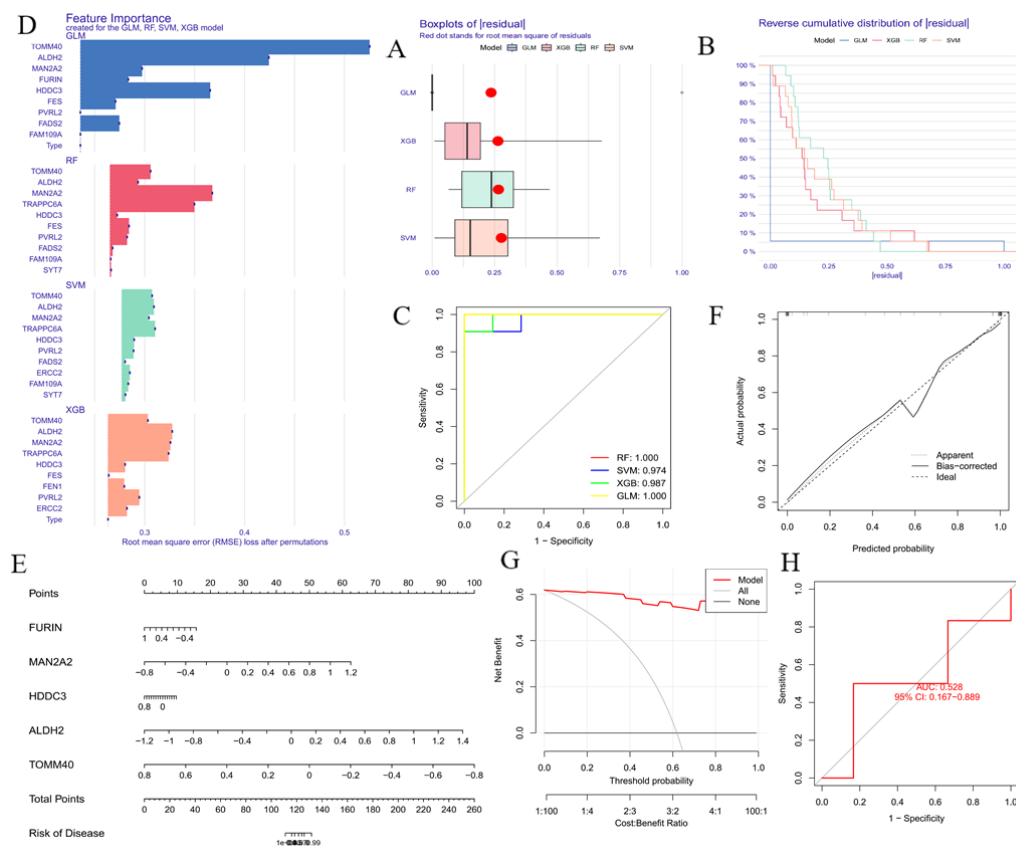
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12

- 1 FIGURE 3 | (A) Bar plot of relative percentage of each immune cells in samples; (B) Box
- 2 plot of immune cell fraction between normal samples and IS samples; (C) Heat map of
- 3 correlation analysis between DEAGs and immune cells



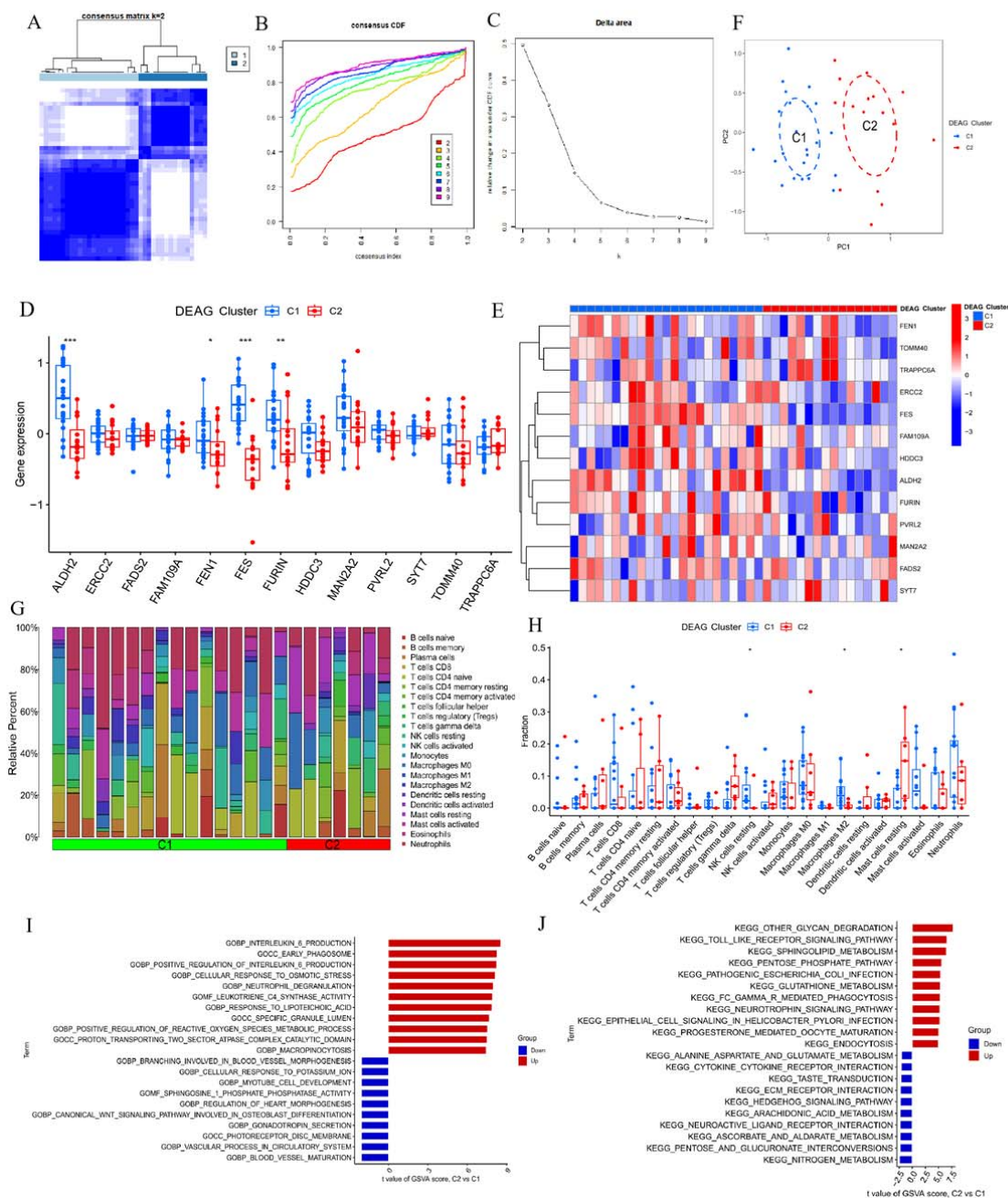
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14

1 FIGURE 4 | (A) Box plots of residual of the four machine learning models; (B) Reverse
 2 cumulative distribution of residual of the four machine learning models; (C) ROC of the four
 3 machine learning models; (D) Bar plot of feature importance of the four machine learning
 4 models; (E) Nomogram of the feature genes; (F) Calibration curve of feature genes
 5 nomogram; (G) Decision curve of feature genes nomogram; (H) ROC of the test GEO dataset
 6



7
 8
 9
 10
 11
 12

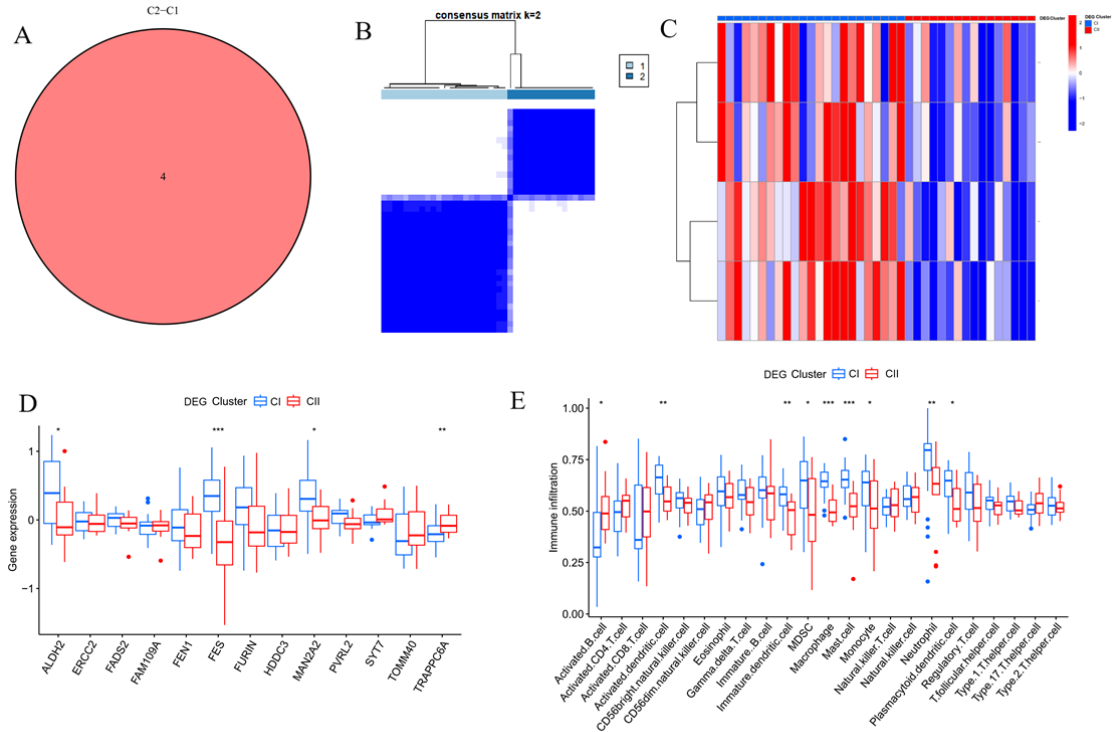
1 FIGURE 5 | (A) Consensus cumulative distribution plot of DEAG clustering for samples; (B)
2 Consensus CDFs of DEAG clustering; (C) Consensus matrix heat map of DEAG clustering
3 for samples; (D) Heat map of DEAGs expression between DEAG clusters; (E) Box plot of
4 expression difference analysis of DEAG clusters; (F) Scatter plot of PCA between DEAG
5 clusters; (G) Box plot of immune cell fraction between DEAG clusters; (H) Bar plot of
6 relative percentage of each immune cells in samples of DEAG clusters; (I) Bar plot of GO
7 terms of GSVA between DEAG clusters; (J) Bar plot of KEGG terms of GSVA between
8 DEAG clusters
9
10
11



1
2
3
4
5
6

1 FIGURE 6 | (A) Venn plot of DEGs; (B) Consensus matrix heat map of DEG clustering for IS
2 samples; (C) Heat map of DEGs expression between DEG clusters; (D) Box plot of
3 expression difference analysis of DEG clusters; (E) Box plot of immune cell infiltration
4 between DEG clusters

5



6

7

8

9

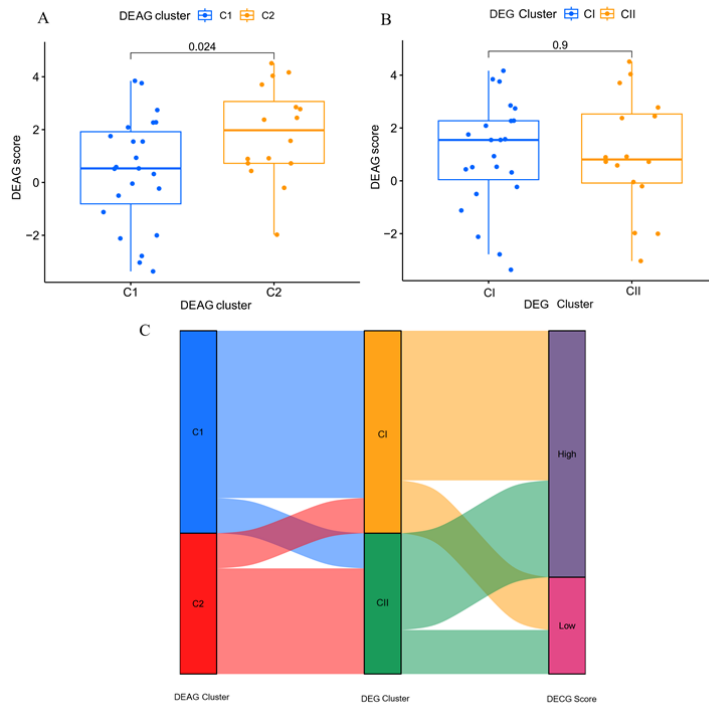
10

11

12

13

- 1 FIGURE 7 | (A) Box plot of different expression analysis of DEAG score between DEAG
 - 2 clusters; (B) Box plot of different expression analysis of DEAG score between DEG clusters;
 - 3 (C) Alluvial plot of the correspondence of the different clustered samples
- 4



5

1 Exploring Genetic Associations of Three Types of Risk Factors with Ischemic Stroke: An
2 Integrated Bioinformatics Study

3

4 Yi Liu^{1#}, PhD; Weili Wang^{1#}, MMedSci; Xin Cui^{1*}, PhD; Yanming Xie^{1*}, Prof.

5

6 ¹Institute of Basic Research in Clinical Medicine, China Academy of Chinese Medical
7 Sciences Beijing, 100700, China;

8

9 #Yi Liu and Weili Wang contributed equally to this work.

10

11 *Yanming Xie and Xin Cui are the joint corresponding authors of this article.

12

13 Correspondence:

14 Yanming Xie No. 16, Nanxiaojie, Dongzhimennei, Dongcheng District, Beijing, China

15 E-mail: ketizu2018@163.com

16 Xin Cui No. 16, Nanxiaojie, Dongzhimennei, Dongcheng District, Beijing, China

17 E-mail: Xinrobertcm@hotmail.com

18

19 Short Title: Associations of exposures and ischemic stroke

20

21 The total number of words: 8817 words

22

23

1 Abstract:

2 Background: Ischemic stroke (IS) is a primary cause of disability and mortality globally.

3 More and more reports suggest a strong association between blood pressure, blood glucose,
4 and blood lipids and their metabolic products with IS.

5 Methods: We extracted the genetic tools of blood pressure, blood glucose, and blood lipids
6 and their metabolites as instrumental variables, which were then paired with GWAS data on
7 IS and a Mendelian randomization (MR) analysis was performed to assess the effect of these
8 exposures on the disease. Following the positive results, colocalization analysis was
9 performed to identify shared genes associated with exposures and IS. We then performed
10 differential expression analysis using the GEO dataset to identify the differentially expressed
11 associated genes (DEAGs) from associated shared genes. Additional analyzes were
12 performed on these DEAGs to obtain their importance scores using four machine learning
13 models. A nomogram was created using genes with high importance scores to predict the
14 level of risk assessment between DEAGs and IS.

15 Results: There is a positive correlation between blood pressure, blood glucose and the risk of
16 IS onset, while blood lipids and their metabolic products are positively or negatively
17 correlated with the risk. There are 64 shared genes of blood pressure, blood lipids and their
18 metabolic products with IS. Thirteen DEAGs were obtained, and among which FURIN,
19 MAN2A2, HDDC3, ALDH2, and TOMM40 were identified as feature genes for creating the
20 nomogram which can quantitatively predict the risk of IS onset with the expression of these
21 feature genes. By cluster analysis, we found that DEAGs expression underlying immune
22 inflammation, angiogenesis and development, lipid metabolism, etc.

1 Conclusion: This study suggests a significant association between blood pressure, blood
2 glucose, and blood lipids and their metabolic products with IS, and predicts that these
3 exposures mainly regulate the occurrence, development, and prognosis of IS through
4 mechanisms such as DNA repair, DNA methylation, mitochondrial repair, apoptosis,
5 autophagy, etc.

6
7 Key words: Mendelian Randomization; Colocalization analysis; Bioinformatics; Ischemic
8 stroke; Blood pressure; Blood glucose; blood lipids and their metabolic products;

9
10 Abbreviations: Area Under the Curve, AUC; Deoxyribonucleic acid, DNA; Differentially
11 expressed associated genes, DEAGs; differentially expressed genes, DEGs; Extreme Gradient
12 Boosting model, XGB; false discovery rate, FDR; Gene Expression Omnibus , GEO; Gene
13 ontology, GO; gene set variation analysis, GSVA; Generalized Linear, GL; genome-wide
14 association study, GWAS; high-density lipoprotein, HDL; Instrumental variables, IVs;
15 Integrated Epidemiology Unit, IEU; intermediate-density lipoprotein, IDL; inverse variance
16 weighted, IVW; Ischemic stroke, IS; Kyoto Encyclopedia of Genes and Genomes, KEGG;
17 low-density lipoprotein, LDL; Mendelian randomization, MR; Natural killer, NK; one-way
18 analysis of variance, ANOVA; posterior probability, PP; Principal Component Analysis, PCA;
19 Random Forest model, RF; receiver operating characteristic curves, ROC; Single nucleotide
20 polymorphisms, SNPs; Single-sample gene set enrichment analysis, ssGSEA; Support Vector
21 Machine model, SVM; very low-density lipoprotein, VLDL

22

23

1 Ischemic stroke (IS) is a disease caused by the blockage of blood vessels in the brain,
2 leading to local cerebral hypoxia and ischemia, resulting in the death of brain cells. This
3 blockage is usually caused by thrombosis or embolism, causing neurological deficits in the
4 brain area with insufficient blood supply.¹⁻³ In 2019, the number of patients who died from
5 stroke worldwide reached 6.55 million, making it the main cause of adult disability and death
6 worldwide.^{4,5} Among newly diagnosed stroke cases, IS accounted for 62.4%.⁶ Therefore,
7 identifying the risk factors for IS is crucial for preventing this disease. Blood pressure, blood
8 glucose, and blood lipids and their metabolic products have received widespread attention as
9 potential factors affecting the risk of IS. A large number of clinical and epidemiological
10 studies have shown that controlling blood pressure, blood glucose, and blood lipids can
11 reduce the risk of IS.⁷⁻⁹ In recent years, with a deeper understanding of the human genome,
12 we have come to realize that the expression of these biomarkers is largely influenced by
13 genetic information.¹⁰ Therefore, exploring high-risk groups from a genetic perspective and
14 identifying shared genes related to blood pressure, blood glucose, and blood lipids and their
15 metabolic products with IS is of great significance for the precise prevention and treatment of
16 IS.

17 Mendelian randomization (MR) is a genetic epidemiological research tool that is widely
18 used to assess the potential causal relationship between exposure and disease.^{11,12} The
19 working mechanism of MR analysis is similar to that of a natural randomized controlled trial.
20 It uses genetic instruments (single nucleotide polymorphisms, SNPs) that are strongly
21 correlated with exposure factors and are not affected by confounding factors under the
22 principle of MR distribution as instrumental variables (IVs).^{13,14} This study used a

1 two-sample MR analysis to explore the potential causal effects of blood pressure, blood
2 glucose, and blood lipids and their metabolic products on IS. The selected datasets are all
3 from the Integrated Epidemiology Unit (IEU) database (<https://gwas.mrcieu.ac.uk/>). This
4 database is public and contains nearly 2.45 trillion genetic associations from 42,334 summary
5 datasets of genome-wide association study (GWAS).

6 Colocalization analysis is a method for assessing whether two correlated traits are driven
7 by the same genetic mechanism. This method often uses the Bayesian model coloc.¹⁵ The
8 coloc model assumes that in each tested region, each trait has at most one association point,
9 and calculates the posterior probability (PP) of all possible association patterns of the two
10 traits through approximate Bayesian factors: 1) H0: no association; 2) H1: only related to trait
11 1; 3) H2: only related to trait 2; 4) H3: related to both trait 1 and trait 2, but through two
12 independent SNPs; 5) H4: related to both trait 1 and trait 2, through a shared SNP. The
13 posterior probabilities of each association pattern are denoted as PP0 to PP4.¹⁵ A higher PP4
14 value (e.g., PP4 > 50%) provides colocalization support, indicating the existence of shared
15 genetic variation between the two traits.¹⁶ This study used colocalization analysis to find
16 shared genes associated with blood pressure, blood glucose, blood lipids and their metabolic
17 products, and IS. These shared genes may provide important clues for revealing the
18 pathogenesis of IS.

19 The Gene Expression Omnibus (GEO) database collects genomic datasets from different
20 species and tissues, which can be used in conjunction with bioinformatics methods to
21 discover evidence of specific gene expression, identify disease prediction factors, and
22 ultimately help us understand the mechanisms by which the genome regulates physiological

1 and pathological states.¹⁷⁻¹⁹ Therefore, based on the associated shared genes of blood pressure,
2 blood glucose, blood lipids and their metabolic products with IS, this study performed
3 differential analysis on normal samples and IS samples in the GEO dataset, and constructed a
4 machine learning model to screen out feature genes among the associated shared genes. On
5 the one hand, it verifies the MR results through real human samples, and on the other hand, it
6 further screens out biomarkers with higher importance, to further explore the mechanisms of
7 blood pressure, blood glucose, blood lipids and their metabolic products on IS.

8

9 1. Study design

10 In this study, aligned with our research objective, we selected GWAS summary datasets
11 for blood pressure, blood glucose, and blood lipids and their metabolic products with IS. We
12 conducted a two-sample MR analysis and chose positive results for further gene
13 colocalization analysis to identify genes associated with the exposure factors and IS.
14 Subsequently, we retrieved two datasets of IS (training and validation sets), encompassing
15 normal and IS groups, from the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>).
16 Differential analysis, correlation analysis, and immune infiltration analysis were carried out
17 on the normal and IS groups of the training set. With the help of machine learning, we
18 derived feature genes related to blood pressure, blood glucose, blood lipids and their
19 metabolic products, and IS. Following this, the results of the machine learning model were
20 validated using the validation set, and two rounds of clustering analysis on the IS group
21 samples were performed to further investigate the functional mechanisms of the
22 exposure-related genes (Figure 1).

1

2 1.1 Data sources

3 In accordance with the goals of this study, we utilized the IEU database, conducting
4 searches with the keywords of "VLDL (very low-density lipoprotein)", "LDL (low-density
5 lipoprotein)", "IDL (intermediate-density lipoprotein)", "HDL (high-density lipoprotein)",
6 "apolipoprotein", "triglyceride", "fatty acid class", "systolic pressure", "diastolic pressure",
7 "blood glucose", "glycosylated hemoglobin", and "IS". This can identify genetic instruments
8 for blood pressure, blood glucose, and blood lipids and their metabolic products as exposure
9 variables. Concurrently, IS-related genetic instruments were selected as outcome variables for
10 an MR analysis, aiming to explore potential causal relationships between exposures and
11 outcomes. By using "ischemic stroke" as the keyword, setting the data type as array-based
12 expression profiles, and specifying the species as homo sapiens, we retrieved samples from
13 the GEO database (<https://www.ncbi.nlm.nih.gov/geo/>) to obtain gene expression and clinical
14 data of IS patients and healthy individuals. Perl code was employed for gene symbol
15 annotation and data adjustment, thus deriving expression levels of genes related to blood
16 pressure, blood glucose, blood lipids and their metabolites, and IS, in both the normal and IS
17 groups.

18

19 1.2 Mendelian Randomization analysis

20 We utilized the "TwoSampleMR" package to conduct a two-sample MR analysis on the
21 relationship between blood pressure, blood glucose, and blood lipids and their metabolic
22 products with IS.²⁰ In the analysis, we set the parameters as follows: 1) The selected genetic

1 instruments need to have a strong correlation with the exposure factors, and the
2 corresponding P-value should be less than 5×10^{-8} ; 2) In conducting Linkage Disequilibrium
3 (LD) clumping, we set the r^2 threshold at 0.01; 3) When performing clumping analysis, we set
4 the window size as 10,000 kilobase pairs.²¹ These stringent parameter settings help us to
5 accurately reveal the causal relationship between exposures and outcomes.

6

7 1.3 Colocalization analysis

8 After completing the two-sample MR analysis, we selected positive results that
9 demonstrated significant associations between the exposure factors and outcome in the MR
10 analysis. We used the "ieugwasr" and "coloc" packages to conduct colocalization analysis on
11 these positive results to determine if the exposure factors and outcome might be associated
12 within the same gene region. We set a posterior probability (PP4) greater than 50% as the
13 standard, indicating that the genes in that region are associated with both exposure and
14 outcomes¹⁶. Finally, we employed the "gassocplot2" package to visualize the colocalization
15 results, clearly demonstrating the associated shared genes of blood pressure, blood glucose,
16 and blood lipids and their metabolic products with IS.

17

18 1.4 Identification of differently expressed associated genes (DEAGs) and analyses of DEAGs

19 Expression levels of genes associated with blood pressure, blood glucose, blood lipids
20 and their metabolic products, and IS were extracted from both the normal and IS groups.
21 Differential expression analysis was performed using the "limma", "pheatmap", "ggpubr" and
22 other R packages, and results were presented as box plots and heatmaps. Genes with a
23 p-value less than 0.05 were defined as differentially expressed associated genes (DEAGs).

1 Using Perl coding, DEAGs were located on chromosomes and their positions were displayed
2 on a circos plot created with the "Rcirco" package. In addition, the "cor" command was used
3 to calculate the correlation coefficient between each two DEAGs and the results were
4 visualized to illustrate the correlation among DEAGs.

5

6 1.5 Analysis of immune cells in IS samples

7 We performed 1000 simulations using the CIBERSORT command in R to obtain a total
8 relative amount of immune cells equal to 1, then visualized the content of immune cells in
9 each sample using a bar plot. Single-sample gene set enrichment analysis (ssGSEA) was
10 performed using the "GSVA" and "GSABase" packages to compare the differences in
11 immune cell contents between the normal group and the IS group, and the results of ssGSEA
12 are presented as box plot. We matched the differentially expressed associated genes (DEAGs)
13 with the ssGSEA scores, performed a correlation test to obtain the correlation coefficient, and
14 then visualized the results as a heat map.

15

16 1.6 Identification of feature genes and construction and validation of nomogram based on ML

17 Expression data from DEAGs was utilized to construct four predictive models: Random
18 Forest model (RF), Support Vector Machine model (SVM), Generalized Linear Model (GL),
19 and Extreme Gradient Boosting model (XGB). The results from these four models were
20 calculated based on the prediction function. Feature genes within DEAGs were then selected
21 through comprehensive analysis of reverse cumulative distribution plots, residual box plots,
22 and receiver operating characteristic (ROC) curves. After selecting the optimal model, line
23 plots were constructed using the expression levels of the feature genes in both the normal and

1 IS groups. Lastly, decision and calibration curves were drawn to assess the accuracy of the
2 line plots. Another dataset that includes both a normal and IS group was obtained from the
3 GEO database, and a machine learning model was constructed using the same R language
4 approach as before. The ROC was plotted to validate the machine learning model constructed
5 in the test dataset.

6

7 1.7 Clustering of DEAGs and analysis between DEAG clusters.

8 The "ConsensusClusterPlus" package in R was used for clustering, employing Euclidean
9 distance type and allowing up to nine clusters. Expression levels between clusters were
10 compared using heatmaps and box plots, and principal component analysis (PCA) was
11 applied to assess differences between clusters. Following this, a ssGSEA was performed on
12 the DEAG clusters, generating bar plots of the amount of individual immune cells in each
13 sample within the different clusters and comparing the differences in immune cell content
14 among different clusters. Gene ontology (GO) and Kyoto Encyclopedia of Genes and
15 Genomes (KEGG) enrichment analyses were performed using gmt files downloaded from the
16 GSEA platform (<http://www.gsea-msigdb.org/>), and gene set variation analysis (GSVA) was
17 conducted in R language to analyze the expression of enrichment items between clusters.
18 Finally, under the filtering conditions of $|\logFC| > 1$ and an adjusted P-value < 0.05 ,
19 differential expression analysis was performed on the gene expression of DEAG clusters.
20 Intersection of DEAG clusters through a Venn diagram yielded differentially expressed genes
21 (DEGs).

22

23 1.8 Clustering of DEGs and analysis between DEG clusters

1 The same clustering method as in section 1.7 was applied to cluster DEGs, and the
2 DEG cluster with the highest precision was selected. Based on the DEGs clustering results,
3 we compared the expression levels of DEGs, the differences in DEAGs expression, and the
4 immune cell content in different clusters. These results were visualized using heatmaps and
5 box plots.

6

7 1.9 Construction of DEAG scores

8 The PCA method was employed to calculate the expression levels of DEAGs for each
9 sample, yielding DEAG scores (The formula is shown below).²² Using R packages such as
10 "limma" and "ggpubr", we performed differential analysis on the scores of DEAGs in both
11 significantly differentially expressed core gene clusters and DEG clusters. We created box
12 plots to illustrate the scores of significantly differentially expressed core genes in samples
13 within the significantly differentially expressed core gene and DEG clusters. In addition, an
14 alluvial diagram was drawn using the package "ggalluvial" to visualize the relationships and
15 overall processes among DEAG clusters, DEG clusters, samples with higher scores, and
16 samples with lower scores of DEAGs.

17

$$\text{DEAG Score} = \sum (PC1_i + PC2_i)$$

18

19 1.10 Statistical analysis

20 This study carried out MR analysis and colocalization analysis using R V4.2.1. In the
21 MR analysis, if the number of SNPs for the exposure factor is less than 5 after filtering, we
22 will use the Wald ratio method for analysis; otherwise, we will opt for the inverse variance

1 weighted (IVW) method for analysis. In colocalization analysis, to maximize the acquisition
2 of shared genes between "exposure" and "outcome", we set the chromosomal locus parameter
3 range from 75 kilobase pairs to 500 kilobase pairs, aiming to include a wider area around the
4 SNP to enhance the likelihood of discovering shared genes related to both exposure and
5 outcome. In the bioinformatics part, we used Strawberry Perl 5.32.1.1 for GEO dataset
6 extraction and data annotation, and all other statistical analyses were performed with R
7 V4.2.1. For two independent samples, we applied a t-test, while for two paired samples we
8 used the Wilcoxon paired-rank sum test, and for three or more groups of data, we employed
9 one-way analysis of variance (ANOVA) and the Kruskal-Wallis rank sum test. The Spearman
10 rank correlation test was used for correlation analysis. We set a P value < 0.05 or a false
11 discovery rate (FDR) < 0.05 after correction by the Benjamini-Hochberg method as being
12 statistically significant.

13

14 2 Results

15

16 2.1 Results of collection of GWAS and GEO datasets

17 We retrieved 516 relevant exposure datasets from the ieu database based on the set
18 keywords, and the IS dataset came from 440,328 participants of European ancestry, including
19 34,217 cases and 406,111 controls. These datasets were used in two-sample MR analysis for
20 IS to explore their potential causal relationships. By searching the GEO database for IS
21 datasets, we filtered out two datasets that met the requirements, namely GSE16561 and
22 GSE22255, using the former as the test set and the latter as the validation set. The former is

1 total RNA data from human peripheral blood samples, including 39 IS samples (diagnosed by
2 MRI) and 24 normal samples; the latter is RNA expression data from human peripheral blood
3 monocytes, including 20 IS samples and 20 normal samples.
4

5 2.2 Results of Mendelian randomization analysis

6 Among the 516 exposure datasets, 114 showed a positive association with IS. These
7 positive results further suggested that increased levels of Apolipoprotein A-I and HDL may
8 reduce the risk of IS, while elevated levels of VLDL, LDL, IDL, Apolipoprotein B,
9 triglycerides, fatty acids, systolic blood pressure, diastolic blood pressure, fasting blood sugar,
10 and glycated hemoglobin might increase the risk of IS. (Please see Attachment 1 for detailed
11 analysis results)
12

13 2.3 Results of colocalization analysis

14 Upon completion of the two-sample MR analysis, we further subjected the 114 positive
15 results to colocalization analysis to study the potential shared genes between the exposure
16 factors and IS. The results showed that the PP4 was greater than 50% for 70 of the positive
17 results, including exposure factors such as LDL, IDL, VLDL, Apolipoprotein B, fatty acids,
18 systolic and diastolic blood pressure. This suggests that these exposure factors may have
19 genes shared with IS. The colocalization results were visualized using the "gassocplot2"
20 package, and we finally identified 64 shared genes between these exposure factors and IS.
21 These genes may play a key role in the association between the selected exposure factors and

1 IS. (Please see Attachment 1 for detailed analysis results and Attachment 2 for specific shared
2 genes situations)

3

4 2.4 Results of DEAGs identification and analysis of DEAGs.

5 Through colocalization analysis, we identified 64 associated shared genes between
6 blood pressure, blood sugar, and blood lipids and their metabolites with IS. Differential
7 analysis between the IS group and the normal group in the test dataset revealed that 13
8 associated shared genes exhibited statistically significant differences. Among them, ERCC2,
9 FEN1, HDDC3, TOMM40, TRAPPC6A were highly expressed in the normal group, while
10 ALDH2, FADS2, FES, FURIN, MAN2A2, PVRL2, SYT7, and FAM109A were highly
11 expressed in the IS group, as shown in Figures A and B. The specific positions of these
12 significantly differentially expressed association genes on the chromosome are shown in
13 Figure C. The correlation analysis between each pair of DEAGs in the IS samples showed
14 that there is a certain correlation between DEAGs, which are mainly positive correlations, as
15 shown in Figures D and E.

16

17 2.5 Results of Immune cell infiltration analysis, immune cell difference analysis and
18 correlation analysis of IS samples

19 Through immune cell infiltration analysis, we obtained the types and content of immune cells
20 expressed in each sample, as shown in Figure A. ssGSEA analysis revealed (Figure B) that
21 the expression differences of 11 types of immune cells in the control group and IS group were
22 statistically significant. Among them, B cells naive, T cells CD8, T cells CD4 memory

1 activated, T cells follicular helper, Natural killer (NK) cells activated, and Dendritic cells
2 activated were highly expressed in the control group, while T cells gamma delta,
3 Macrophages M0, Macrophages M2, Mast cells activated, and Neutrophils were highly
4 expressed in the IS group. The correlation analysis between DEAGs and immune cells
5 showed that there is a certain correlation between DEAGs and immune cells, and the levels of
6 positive and negative correlations are comparable, as shown in Figure C.

7

8 2.6 Results of selection of machine learning models, construction of nomogram, and
9 verification

10 The results of building SVM, RF, XGB, and GL machine learning prediction models
11 using DEAGs data were obtained. From the residual boxplots, inverse cumulative distribution
12 plots, and ROC curves (Figure A, C, and D), it can be observed that the GL method had the
13 highest area under the ROC curve, the lowest residual values, and the lowest inverse
14 cumulative values. Therefore, the GL method was considered the most accurate and chosen
15 as the best model for further analysis. The GL model provided importance scores for the
16 selected feature genes, as shown in Figure B, revealing nine feature genes. Among them, the
17 top five genes with the highest importance scores were used to construct the column line
18 graph (including FURIN, MAN2A2, HDDC3, ALDH2, and TOMM40). Then, separate
19 scoring scales were obtained for these five feature genes. The risk rate of the co-located
20 feature genes associated with blood pressure, blood glucose, and blood lipids and their
21 metabolic products in the occurrence of IS was assessed by calculating the sum of the feature
22 gene expression scores (Figure E). The predictive accuracy of the model was evaluated by the

1 distance between the solid line and the dotted line in the calibration curve (Figure F) and the
2 distance between the red and gray lines in the decision curve (Figure G), indicating high
3 accuracy. The GL model was validated using the validation dataset (GSE22255) and the top
4 five feature genes with the highest importance scores were included in the model validation
5 and ROC curve analysis. The results showed an Area Under the Curve (AUC) value of 0.528,
6 with a 95% confidence interval of 0.167-0.889. Therefore, based on the validation of the
7 model, it can be concluded that the model based on the GEO dataset has good accuracy.

8

9 2.7 Results of Clustering of DEAGs, Analysis of DEAGs Expression between Clusters,
10 Immune Cell Analysis, and GO and KEGG Enrichment Analysis

11 Clustering analysis based on the expression of DEAGs revealed two distinct clusters
12 with the highest accuracy, as shown in Figure A. Consequently, the IS samples were divided
13 into two groups: C1 and C2, as depicted in Figures B and C. Subsequently, DEAGs
14 expression analysis was performed between the two DEAG clusters (Figure D and E),
15 indicating that ALDH2, FEN1, FES, and FURIN were significantly upregulated in C1. PCA
16 analysis (Figure F) demonstrated that DEAGs can differentiate between C1 and C2.
17 Moreover, ssGSEA analysis (Figure G) identified three statistically significant immune cell
18 types between C1 and C2. Specifically, NK cells resting and Macrophages M2 were
19 upregulated in C1, while Mast cells resting were upregulated in C2. The respective content of
20 different immune cells in each C1 and C2 sample is illustrated in Figure H.

21 In GSVA analysis (Figures I and J), it can be observed that compared to C1, C2 showed
22 upregulation of several GO biological processes, including Interleukin 6 production, positive

1 regulation of interleukin 6 production, cellular response to osmotic stress, neutrophil
2 degeneration, response to lipoteichoic acid, positive regulation of reactive oxygen species
3 metabolic process, and macropinocytosis. On the other hand, C2 exhibited downregulation of
4 branching involved in blood vessel morphogenesis, cellular response to potassium ion,
5 myotube cell development, regulation of heart morphogenesis, canonical Wnt signaling
6 pathway involved in osteoblast differentiation, gonadotropin secretion, vascular process in
7 the circulatory system, and blood vessel maturation in GO biological processes. In GO
8 molecular functions, C2 showed upregulation of leukotriene C4 synthase activity and
9 downregulation of sphingosine 1-phosphate phosphatase activity. In GO cellular components,
10 C2 exhibited upregulation of early phagosome, specific granule lumen, and
11 proton-transporting two-sector ATPase complex catalytic domain, while downregulation of
12 photoreceptor disc membrane. In KEGG pathways, C2 demonstrated upregulation of other
13 glycan degradation, toll-like receptor signaling pathway, sphingolipid metabolism, pentose
14 phosphate pathway, pathogenic *Escherichia coli* infection, glutathione metabolism, Fc gamma
15 R-mediated phagocytosis, neurotrophin signaling pathway, cell signaling in *Helicobacter*
16 *pylori* infection, progesterone-mediated oocyte maturation, and endocytosis. Conversely, C2
17 showed downregulation of alanine, aspartate, and glutamate metabolism, cytokine-cytokine
18 receptor interaction, taste transduction, ECM receptor interaction, hedgehog signaling
19 pathway, arachidonic acid metabolism, neuroactive ligand-receptor interaction, ascorbate and
20 aldarate metabolism, pentose and glucuronate interconversions, and nitrogen metabolism in
21 KEGG pathways.

22

1 2.8 Results of DEGs screening, DEGs clustering, and analyzes of DEG clusters

2 After filtering for significantly differentially expressed genes (DEGs) between C1 and
3 C2 samples, a total of 4 DEGs were identified (Figure A). These 4 DEGs were further
4 subjected to clustering analysis within the DEAGs cluster, resulting in two distinct clusters,
5 namely CI and CII (Figure B). Analysis of DEAGs expression in these two clusters in IS
6 samples revealed that RNASE2, RNASE3, MMP9, and CAMP exhibited statistically
7 significant differential expression, with higher expression in CI and lower expression in CII
8 (Figure C).

9 The differential expression analysis of DEAGs based on DEGs clustering was
10 performed on the samples, as shown in Figure D. The results revealed that these genes were
11 mainly upregulated in the CI group and downregulated in the CII group. Among them, the
12 significantly differentially expressed genes included ALDH2, FES, MAN2A2 (upregulated in
13 CI), and TRAPPCA (upregulated in CII). Additionally, using ssGSEA based on DEGs
14 clustering (using the same gene set file as above), statistically significant differences in
15 immune cells were observed, including activated B cells (downregulated in CI), activated
16 dendritic cells, immature dendritic cells, myeloid-derived suppressor cells (MDSCs),
17 macrophages, mast cells, monocytes, neutrophils, and plasmacytoid dendritic cells
18 (upregulated in CI), as shown in Figure E.

19

20 2.9 Results of DEAGs scoring, differential analysis of DEAGs score between clusters, and 21 construction of the alluvial plot

22 The differential analysis of DEGs clustering based on the principal component analysis

1 (PCA) scores showed statistically significant differences between the clusters (Figure A).
2 Specifically, C1 had lower scores, while C2 had higher scores. However, there was no
3 statistically significant difference observed in the DEGs clustering (Figure B). The scatter
4 plot in Figure C displayed that the C1 and C2 clusters of DEAGs corresponded mainly to the
5 CI and CII clusters of DEGs, respectively. However, there was no clear correspondence
6 observed between the high and low scores of DEAGs and the DEGs clustering.

7

8 3 Discussion

9

10 3.1. The relationships of blood pressure, blood glucose, and blood lipids and their metabolic
11 products with IS

12 Clinical studies have demonstrated that effective management of various physiological
13 indicators can significantly impact the risk of IS. For example, blood pressure control has
14 been shown to significantly reduce the risk of IS recurrence.²³ Elevated fasting blood glucose
15 levels are significantly associated with an increased risk of IS.²⁴ Furthermore, increased
16 levels of glycated hemoglobin, regardless of diabetes status, are associated with an increased
17 risk of IS.²⁵ Lipid management also plays a significant role in IS risk. There are significant
18 differences in VLDL levels between the control and IS groups.²⁶ However, the role and
19 mechanisms of VLDL level changes in IS pathogenesis require further investigation.
20 Medications that lower LDL levels can significantly reduce the risk of IS.²⁷ Increased levels
21 of IDL and apolipoprotein B are associated with an increased risk of IS.^{28,29} Controlling
22 triglyceride levels can lower the risk of IS.³⁰ Both polyunsaturated fatty acids and saturated

1 fatty acids are also associated with IS.^{31,32} However, elevated levels of HDL and
2 apolipoprotein A-I are associated with a lower risk of IS.^{33,34} Taking all the above clinical
3 studies into consideration, it can be concluded that blood pressure, blood glucose, and blood
4 lipids and their metabolic products play critical roles in influencing the risk of IS.

5 The results of MR analysis indicate a potential causal relationship between blood
6 pressure (diastolic and systolic), blood glucose (fasting blood glucose and glycated
7 hemoglobin), blood lipids and their various metabolic products (VLDL, LDL, IDL, HDL,
8 apolipoprotein A-I, apolipoprotein B, triglycerides, Eicosapentaenoate, Docosapentaenoate,
9 Stearidonate, Docosahexaenoic acid) with IS. Increasing levels of high-density lipoprotein
10 (HDL) and apolipoprotein A-I are associated with a decreased risk of IS. However, elevated
11 blood pressure (including systolic and diastolic), blood glucose (fasting blood glucose and
12 glycated hemoglobin), various low-density lipoproteins such as VLDL, LDL, IDL, as well as
13 apolipoprotein B, triglycerides, Eicosapentaenoate, Docosapentaenoate, Stearidonate,
14 Docosahexaenoic acid are associated with an increased risk of IS.

15

16 3.2 Discussion of mechanisms associated shared genes

17 Through MR analysis, we found that blood pressure, blood glucose, blood lipids and
18 their metabolic products have a significant impact on the risk of IS. To further explore the
19 specific mechanisms underlying these effects, we conducted a colocalization analysis. The
20 results showed that there are 64 associated shared genes between LDL, IDL, VLDL,
21 apolipoprotein B, fatty acids, systolic blood pressure, and diastolic blood pressure with IS.
22 Although a positive association between blood glucose and IS risk has been established, we

1 did not identify any associated shared genes due to the PP4 being less than 50%.

2 The differential analysis of the 64 associated shared genes between normal and IS
3 samples, using the GEO dataset, revealed 13 genes with significant differential expression.
4 Among them, ERCC2, a transcription factor (TF) involved in nucleotide excision repair of
5 damaged DNA,³⁵ has been clinically associated with increased stroke risk,^{36,37} and animal
6 studies have shown its neuroprotective role in preventing ischemia-reperfusion injury.³⁸
7 FEN1, another TF encoding a protein involved in rDNA and mitochondrial DNA repair,³⁹ has
8 been implicated in the pathogenesis of IS and is considered a therapeutic target for IS
9 treatment.^{40,41} HDDC3, a protein encoded by HDDC3 gene, participates in the starvation
10 response, and starvation-induced autophagy is known to protect neurons and regulate IS.^{42,43}
11 TOMM40, encoding the translocase of the outer mitochondrial membrane 40 homolog, has
12 been identified as a potential susceptibility locus for IS based on a study conducted in
13 Japan.⁴⁴ ALDH2, encoding aldehyde dehydrogenase, has shown neuroprotective effects by
14 clearing 4-hydroxy-2-nonenal and reducing mitochondrial-associated cell apoptosis through
15 JNK-mediated cystathionine- β -synthase-3 activation, making it a potential target for IS
16 intervention. Case-control studies have also suggested an association between ALDH2
17 polymorphisms and IS risk in the Han Chinese population.⁴⁵⁻⁴⁷ FADS2, encoding fatty acid
18 desaturase 2, has been associated with IS risk and lipid levels in case-control studies,
19 although the specific mechanisms are still under investigation.⁴⁸ The product of the FES gene
20 exhibits tyrosine-specific protein kinase activity. Inhibitors of this class of enzymes have
21 shown potential therapeutic value for various diseases, including IS, but the involvement of
22 FES in IS is yet to be confirmed.⁴⁹ FURIN, encoding proprotein convertase subtilisin/kexin

1 type 1, has been found to be upregulated within 24 hours after cerebral ischemia in animal
2 experiments,⁵⁰ and case-control studies have revealed lower DNA methylation levels in the IS
3 region, suggesting its association with IS.⁵¹ SYT7, encoding synaptotagmin 7, has been
4 identified to have significant DNA methylation correlation with IS in a large-scale
5 sequencing study.⁵² The roles of FAM109A, MAN2A2, PVRL2, and others in IS require
6 further investigation.

7 These genes primarily function through DNA repair, DNA methylation, mitochondrial
8 repair, cell apoptosis, and autophagy to regulate lipid levels, blood pressure, and neuronal
9 protection, thereby affecting the occurrence, development, and prognosis of IS. The
10 associated shared genes are predominantly upregulated in the IS group, and they exhibit
11 mostly positive regulatory relationships, indicating potential synergistic effects among these
12 genes in IS. Differential expression analysis of immune cells between the IS and normal
13 groups reveals significant expression differences in 50% of the immune cells, suggesting the
14 involvement of immune cells in the regulation of IS by these associated shared genes.
15 Moreover, the comparable levels of positive and negative regulation of immune cells by
16 DEAGs suggest that DEAGs may have bidirectional regulatory effects on the immune
17 system's role in IS.

18 The feature genes obtained from the GL model constructed through machine learning
19 (FURIN, MAN2A2, HDDC3, ALDH2, TOMM40) have significant importance in regulating
20 the associated shared genes and the risk of developing IS. The column plot provides a
21 quantitative prediction of the aforementioned importance and risk of disease. Furthermore,

1 the construction of the GL model using the validation dataset indicates that the column plot
2 has high accuracy.

3 The IS samples were clustered into C1 and C2 groups based on the expression of
4 DEAGs, with C1 showing predominantly high expression of DEAGs and C2 showing low
5 expression. The two clusters were significantly distinguishable according to the PCA results.
6 Furthermore, immune cell infiltration analysis of DEAGs clustering revealed significant
7 differences in immune cell expression between C1 and C2 groups of IS patients. The GSVA
8 results indicated significant expression differences between C1 and C2 groups in various
9 biological processes and pathways, including immune inflammation, vascular development
10 and formation, reactive oxygen species metabolism, carbohydrate metabolism, lipid
11 metabolism, fatty acid metabolism, amino acid metabolism, and neural regulation.

12 Differential analysis between the C1 and C2 groups identified 4 DEGs, which have been
13 extensively studied and shown to be associated with IS incidence, IS prognosis, and IS
14 complications.⁵³⁻⁵⁵ Subsequently, IS samples were further clustered into CI and CII groups
15 based on the expression of DEGs. Differential analysis revealed that DEGs were highly
16 expressed in the CI group and lowly expressed in the CII group, while DEAGs were mainly
17 upregulated in the CI group and downregulated in the CII group. Immune cell infiltration
18 analysis of DEGs clustering demonstrated significant differences in immune cell expression
19 between the CI and CII groups of IS patients, with upregulation observed in the CI group and
20 downregulation observed in the CII group. PCA-based DEAG scoring and subsequent
21 differential analysis revealed significant differences in DEAG scoring between the C1 and C2
22 groups, while no significant differences were observed between the CI and CII groups.

1 Furthermore, the construction of Sankey diagrams provided a visual representation of the
2 correspondence between the two clustering results and the high/low DEAGs\ scoring groups,
3 with C1 and C2 corresponding to CI, CII, and the high/low DEAG groups, while the
4 correspondence between CI and CII was less clear.

6 3.3 Limitations of this study

7 Although MR analysis and colocalization analysis provide valuable insights into the risk
8 factors and underlying mechanisms of IS from the perspective of exposure factors and genes,
9 these methods also have their limitations. The pathogenesis of IS is extremely complex and
10 involves intricate interactions among multiple environmental factors, physiological indicators,
11 and genetic factors. Therefore, in order to comprehensively understand and elucidate the risk
12 factors and mechanisms of IS, a holistic and multifactorial research approach is needed.
13 Additionally, the GWAS data included in our study primarily originate from sample
14 populations of European ancestry. Hence, extrapolating these results directly to non-European
15 populations may have limitations. Similarly, bioinformatics methods also have certain
16 limitations. Although they have provided multi-faceted and multi-dimensional analyses of
17 human blood samples, elucidating specific mechanisms of interaction between blood pressure,
18 blood glucose, blood lipids and their metabolic products in relation to IS, further research and
19 clinical applications require careful consideration. Moreover, the blood samples used in this
20 study were only collected from the United States and Portugal, and whether the results can be
21 extrapolated to other countries, ethnicities, and populations remains to be investigated and
22 confirmed. Given the relatively small sample size, there is a greater risk of bias in the
23 conclusions. Therefore, large-scale studies or animal experiments can be conducted to further

1 validate and deepen the conclusions.

2

3 4 Conclusion and prospects

4 This study first integrated MR analysis, colocalization analysis, and bioinformatics
5 analysis to conduct an in-depth study of the association of blood pressure, blood glucose, and
6 blood lipids and their metabolic products with IS. Firstly, the results of MR analysis showed
7 that there is a positive correlation between blood pressure, blood glucose and the risk of IS
8 onset, and blood lipids and their metabolic products have either a positive or negative
9 correlation with the risk of IS onset. Furthermore, through colocalization analysis, we
10 identified 64 associated shared genes with blood pressure, blood lipids and their metabolic
11 products in relation to IS. Additionally, using the GEO dataset, we identified 13 DEAGs.
12 These genes primarily function through deoxyribonucleic acid (DNA) repair, DNA
13 methylation, mitochondrial repair, cell apoptosis, and autophagy, regulating lipid metabolism,
14 blood pressure, and neuroprotection, thereby influencing the occurrence, progression, and
15 prognosis of IS. Subsequently, through the construction of machine learning models, we
16 selected feature genes. Cluster analysis revealed that the expression of DEAGs in IS may also
17 involve mechanisms such as immune inflammation, vascular development, lipid metabolism,
18 and fatty acid metabolism.

19 The findings of this study will contribute to a deeper understanding of the
20 pathophysiological mechanisms underlying IS, optimize treatment strategies, and drive the
21 development of new drugs. The methodological approach used in this study, which combines
22 genetic epidemiology with bioinformatics, offers a unique perspective for elucidating causal

1 relationships between exposure factors and diseases, as well as unraveling genetic
2 mechanisms. This research strategy, along with the identification of the associated shared
3 genes, provides a new theoretical framework and practical approach for the prevention,
4 diagnosis, and treatment of IS.

5

6 Article Information

7

8 Acknowledgments

9 The authors would like to thanks the institutions for providing GWAS summary statistics and
10 human gene expression datasets.

11

12 Source of Founding

13 National Administration of Traditional Chinese Medicine 2021 Qihuang Scholar Support
14 Project (National TCM Education Letter 2022 No. 6) ; Yanming Xie's National Renowned
15 Traditional Chinese Medicine Expert Inheritance Studio Construction Project (National TCM
16 Education Letter 2022 No. 75).

17

18 Disclosures

19 The authors declare that the research was conducted in the absence of any commercial or
20 financial relationships that could be construed as a potential conflict of interest.

21

22 References:

- 1 1. Hankey GJ. Stroke. *Lancet*. 2017;389:641-654. doi: 10.1016/S0140-6736(16)30962-X.
- 2 2. Lee RHC, Lee MHH, Wu CYC, Couto E Silva A, Possoit HE, Hsieh TH, Minagar A, Lin
3 HW. Cerebral ischemia and neuroregeneration. *Neural Regen Res*. 2018;13:373-385. doi:
4 10.4103/1673-5374.228711.
- 5 3. Powers WJ, Rabinstein AA, Ackerson T, Adeoye OM, Bambakidis NC, Becker K, Biller J,
6 Brown M, Demaerschalk BM, Hoh B, et al. Guidelines for the Early Management of
7 Patients With Acute Ischemic Stroke: 2019 Update to the 2018 Guidelines for the Early
8 Management of Acute Ischemic Stroke: A Guideline for Healthcare Professionals From the
9 American Heart Association/American Stroke Association. *Stroke*. 2019;50:e344-e418. doi:
10 10.1161/STR.0000000000000211.
- 11 4. Romano JG, Rundek T. Expanding Treatment for Acute Ischemic Stroke beyond
12 Revascularization. *N Engl J Med*. 2023;388:2095-2096. doi: 10.1056/NEJMe2303184.
- 13 5. Roth GA, Mensah GA, Johnson CO, Addolorato G, Ammirati E, Baddour LM, Barengo
14 NC, Beaton AZ, Benjamin EJ, Benziger CP, et al. Global Burden of Cardiovascular
15 Diseases and Risk Factors, 1990-2019: Update From the GBD 2019 Study. *J Am Coll*
16 *Cardiol*. 2020;76:2982-3021. doi: 10.1016/j.jacc.2020.11.010.
- 17 6. GBD 2019 Stroke Collaborators. Global, regional, and national burden of stroke and its
18 risk factors, 1990-2019: a systematic analysis for the Global Burden of Disease Study 2019.
19 *Lancet Neurol*. 2021;20:795-820. doi: 10.1016/S1474-4422(21)00252-0.
- 20 7. Hansson L, Zanchetti A, Carruthers SG, Dahlöf B, Elmfeldt D, Julius S, Ménard J, Rahn
21 KH, Wedel H, Westerling S. Effects of intensive blood-pressure lowering and low-dose
22 aspirin in patients with hypertension: principal results of the Hypertension Optimal

- 1 Treatment (HOT) randomised trial. HOT Study Group. *Lancet*. 1998;351:1755-1762. doi:
2 10.1016/s0140-6736(98)04311-6.
- 3 8. Duckworth W, Abaira C, Moritz T, Reda D, Emanuele N, Reaven PD, Zieve FJ, Marks J,
4 Davis SN, Hayward R, et al. Glucose control and vascular complications in veterans with
5 type 2 diabetes. *N Engl J Med*. 2009;360:129-139. doi: 10.1056/NEJMoa0808431.
- 6 9. Amarenco P, Labreuche J. Lipid management in the prevention of stroke: review and
7 updated meta-analysis of statins for stroke prevention. *Lancet Neurol*. 2009;8:453-463. doi:
8 10.1016/S1474-4422(09)70058-4.
- 9 10. Kim YS, Leventhal BL. Genetic epidemiology and insights into interactive genetic and
10 environmental effects in autism spectrum disorders. *Biol Psychiatry*. 2015;77:66-74. doi:
11 10.1016/j.biopsych.2014.11.001.
- 12 11. Smith GD, Ebrahim S. 'Mendelian randomization': can genetic epidemiology contribute to
13 understanding environmental determinants of disease? *Int J Epidemiol*. 2003;32:1-22. doi:
14 10.1093/ije/dyg070.
- 15 12. Grover S, Del Greco M F, Stein CM, Ziegler A. Mendelian Randomization. *Methods Mol*
16 *Biol*. 2017;1666:581-628. doi: 10.1007/978-1-4939-7274-6_29.
- 17 13. Verduijn M, Siegerink B, Jager KJ, Zoccali C, Dekker FW. Mendelian randomization: use
18 of genetics to enable causal inference in observational studies. *Nephrol Dial Transplant*.
19 2010;25:1394-1398. doi: 10.1093/ndt/gfq098.
- 20 14. Davies NM, Holmes MV, Davey Smith G. Reading Mendelian randomisation studies: a
21 guide, glossary, and checklist for clinicians. *BMJ*. 2018;362:k601. doi: 10.1136/bmj.k601.
- 22 15. Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, Plagnol

- 1 V. Bayesian test for colocalisation between pairs of genetic association studies using
2 summary statistics. *PLoS Genet.* 2014;10:e1004383. doi: 10.1371/journal.pgen.1004383.
- 3 16. Yarmolinsky J, Bouras E, Constantinescu A, Burrows K, Bull CJ, Vincent EE, Martin RM,
4 Dimopoulou O, Lewis SJ, Moreno V, et al. Genetically proxied glucose-lowering drug
5 target perturbation and risk of cancer: a Mendelian randomisation analysis. *Diabetologia.*
6 2023. doi: 10.1007/s00125-023-05925-4.
- 7 17. Barrett T, Troup DB, Wilhite SE, Ledoux P, Rudnev D, Evangelista C, Kim IF, Soboleva
8 A, Tomashevsky M, Marshall KA, et al. NCBI GEO: archive for high-throughput
9 functional genomic data. *Nucleic Acids Res.* 2009;37(Database issue):D885-890. doi:
10 10.1093/nar/gkn764.
- 11 18. Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, Marshall KA,
12 Phillippy KH, Sherman PM, Holko M, et al. NCBI GEO: archive for functional genomics
13 data sets--update. *Nucleic Acids Res.* 2013;41(Database issue):D991-995. doi:
14 10.1093/nar/gks1193.
- 15 19. Xu B, Wang L, Zhan H, Zhao L, Wang Y, Shen M, Xu K, Li L, Luo X, Zhou S, et al.
16 Investigation of the Mechanism of Complement System in Diabetic Nephropathy via
17 Bioinformatics Analysis. *J Diabetes Res.* 2021;2021:5546199. doi: 10.1155/2021/5546199.
- 18 20. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, Laurin C, Burgess S,
19 Bowden J, Langdon R, et al. The MR-Base platform supports systematic causal inference
20 across the human phenome. *Elife.* 2018;7:e34408. doi: 10.7554/eLife.34408.
- 21 21. Xiang M, Wang Y, Gao Z, Wang J, Chen Q, Sun Z, Liang J, Xu J. Exploring causal
22 correlations between inflammatory cytokines and systemic lupus erythematosus: A

- 1 Mendelian randomization. *Front Immunol.* 202;13:985729. doi:
2 10.3389/fimmu.2022.985729.
- 3 22. Zhang B, Wu Q, Li B, Wang D, Wang L, Zhou YL. m6A regulator-mediated methylation
4 modification patterns and tumor microenvironment infiltration characterization in gastric
5 cancer. *Mol Cancer.* 2020;19:53. doi: 10.1186/s12943-020-01170-0.
- 6 23. Biffi A, Anderson CD, Battey TW, Ayres AM, Greenberg SM, Viswanathan A, Rosand J.
7 Association Between Blood Pressure Control and Risk of Recurrent Intracerebral
8 Hemorrhage. *JAMA.* 2015;314:904-912. doi: 10.1001/jama.2015.10082.
- 9 24. Boden-Albala B, Cammack S, Chong J, Wang C, Wright C, Rundek T, Elkind MS, Paik
10 MC, Sacco RL. Diabetes, fasting glucose levels, and risk of ischemic stroke and vascular
11 events: findings from the Northern Manhattan Study (NOMAS). *Diabetes Care.*
12 2008;31:1132-1137. doi: 10.2337/dc07-0797.
- 13 25. Nomani AZ, Nabi S, Ahmed S, Iqbal M, Rajput HM, Rao S. High HbA1c is associated
14 with higher risk of ischaemic stroke in Pakistani population without diabetes. *Stroke Vasc*
15 *Neurol.* 2016;1:133-139. doi: 10.1136/svn-2016-000018.
- 16 26. Bansal BC, Sood AK, Bansal CB. Familial hyperlipidemia in stroke in the young. *Stroke.*
17 1986;17:1142-1145. doi: 10.1161/01.str.17.6.1142.
- 18 27. Schwartz GG, Steg PG, Szarek M, Bhatt DL, Bittner VA, Diaz R, Edelberg JM,
19 Goodman SG, Hanotin C, Harrington RA, et al. Alirocumab and Cardiovascular Outcomes
20 after Acute Coronary Syndrome. *N Engl J Med.* 2018;379:2097-2107. doi:
21 10.1056/NEJMoa1801174.

- 1 28. Berger JS, McGinn AP, Howard BV, Kuller L, Manson JE, Otvos J, Curb JD, Eaton CB,
2 Kaplan RC, Lynch JK, et al. Lipid and lipoprotein biomarkers and the risk of ischemic
3 stroke in postmenopausal women. *Stroke*. 2012;43:958-66. doi:
4 10.1161/STROKEAHA.111.641324.
- 5 29. Johannesen CDL, Mortensen MB, Langsted A, Nordestgaard BG. ApoB and Non-HDL
6 Cholesterol Versus LDL Cholesterol for Ischemic Stroke Risk. *Ann Neurol*.
7 2022;92:379-389. doi: 10.1002/ana.26425.
- 8 30. Das Pradhan A, Glynn RJ, Fruchart JC, MacFadyen JG, Zaharris ES, Everett BM,
9 Campbell SE, Oshima R, Amarenco P, Blom DJ, et al. Triglyceride Lowering with
10 Pemafibrate to Reduce Cardiovascular Risk. *N Engl J Med*. 2022 ;387:1923-1934. doi:
11 10.1056/NEJMoa2210645.
- 12 31. Saber H, Yakoob MY, Shi P, Longstreth WT Jr, Lemaitre RN, Siscovick D, Rexrode KM,
13 Willett WC, Mozaffarian D. Omega-3 Fatty Acids and Incident Ischemic Stroke and Its
14 Atherothrombotic and Cardioembolic Subtypes in 3 US Cohorts. *Stroke*.
15 2017;48:2678-2685. doi: 10.1161/STROKEAHA.117.018235.
- 16 32. Venø SK, Schmidt EB, Jakobsen MU, Lundbye-Christensen S, Bach FW, Overvad K.
17 Substitution of Linoleic Acid for Other Macronutrients and the Risk of Ischemic Stroke.
18 *Stroke*. 2017;48:3190-3195. doi: 10.1161/STROKEAHA.117.017935.
- 19 33. Curb JD, Abbott RD, Rodriguez BL, Masaki KH, Chen R, Popper JS, Petrovitch H, Ross
20 GW, Schatz IJ, Belleau GC, et al. High density lipoprotein cholesterol and the risk of
21 stroke in elderly men: the Honolulu heart program. *Am J Epidemiol*. 2004;160:150-157.
22 doi: 10.1093/aje/kwh177.

- 1 34. O'Donnell MJ, McQueen M, Sniderman A, Pare G, Wang X, Hankey GJ, Rangarajan S,
2 Chin SL, Rao-Melacini P, Ferguson J, et al. Association of Lipids, Lipoproteins, and
3 Apolipoproteins with Stroke Subtypes in an International Case Control Study
4 (INTERSTROKE). *J Stroke*. 2022;24:224-235. doi: 10.5853/jos.2021.02152.
- 5 35. National Center for Biotechnology Information (2023). PubChem Gene Summary for
6 Gene 2068, ERCC2 - ERCC excision repair 2, TFIIH core complex helicase subunit
7 (human). Retrieved June 14, 2023
8 from <https://pubchem.ncbi.nlm.nih.gov/gene/ERCC2/human>.
- 9 36. Shyu HY, Shieh JC, Ji-Ho L, Wang HW, Cheng CW. Polymorphisms of DNA repair
10 pathway genes and cigarette smoking in relation to susceptibility to large artery
11 atherosclerotic stroke among ethnic Chinese in Taiwan. *J Atheroscler Thromb*.
12 2012;19:316-325. doi: 10.5551/jat.10967.
- 13 37. LD, Dawsey SM, Dong ZW, Taylor PR, Mark SD. A prospective study of polymorphisms
14 of DNA repair genes XRCC1, XPD23 and APE/ref-1 and risk of stroke in Linxian, China.
15 *J Epidemiol Community Health*. 2007;61:737-741. doi: 10.1136/jech.2006.048934.
- 16 38. Zhang J, Guo F, Zhou R, Xiang C, Zhang Y, Gao J, Cao G, Yang H. Proteomics and
17 transcriptome reveal the key transcription factors mediating the protection of Panax
18 notoginseng saponins (PNS) against cerebral ischemia/reperfusion injury. *Phytomedicine*.
19 2021;92:153613. doi: 10.1016/j.phymed.2021.153613.
- 20 39. National Center for Biotechnology Information (2023). PubChem Gene Summary for
21 Gene 2237, FEN1 - flap structure-specific endonuclease 1 (human). Retrieved June 15,
22 2023 from <https://pubchem.ncbi.nlm.nih.gov/gene/FEN1/human>.

- 1 40. He Z, Ning N, Zhou Q, Khoshnam SE, Farzaneh M. Mitochondria as a therapeutic target
2 for ischemic stroke. *Free Radic Biol Med.* 2020;146:45-58. doi:
3 10.1016/j.freeradbiomed.2019.11.005.
- 4 41. Andrabi SS, Parvez S, Tabassum H. Ischemic stroke and mitochondria: mechanisms and
5 targets. *Protoplasma.* 2020;257:335-343. doi: 10.1007/s00709-019-01439-2.
- 6 42. Sun D, Lee G, Lee JH, Kim HY, Rhee HW, Park SY, Kim KJ, Kim Y, Kim BY, Hong JI,
7 et al. A metazoan ortholog of SpoT hydrolyzes ppGpp and functions in starvation
8 responses. *Nat Struct Mol Biol.* 2010;17:1188-94. doi: 10.1038/nsmb.1906.
- 9 43. He S, Wang C, Dong H, Xia F, Zhou H, Jiang X, Pei C, Ren H, Li H, Li R, et al.
10 Immune-related GTPase M (IRGM1) regulates neuronal autophagy in a mouse model of
11 stroke. *Autophagy.* 2012;8:1621-1627. doi: 10.4161/auto.21561.
- 12 44. Yamase Y, Horibe H, Ueyama C, Fujimaki T, Oguri M, Kato K, Arai M, Watanabe S,
13 Yamada Y. Association of *TOMM40* and *SLC22A4* polymorphisms with ischemic stroke.
14 *Biomed Rep.* 2015;3:491-498. doi: 10.3892/br.2015.457.
- 15 45. Guo JM, Liu AJ, Zang P, Dong WZ, Ying L, Wang W, Xu P, Song XR, Cai J, Zhang SQ,
16 et al. ALDH2 protects against stroke by clearing 4-HNE. *Cell Res.* 2013;23:915-930. doi:
17 10.1038/cr.2013.69.
- 18 46. Xia P, Zhang F, Yuan Y, Chen C, Huang Y, Li L, Wang E, Guo Q, Ye Z. ALDH 2
19 conferred neuroprotection on cerebral ischemic injury by alleviating mitochondria-related
20 apoptosis through JNK/caspase-3 signing pathway. *Int J Biol Sci.* 2020;16:1303-1323. doi:
21 10.7150/ijbs.38962.
- 22 47. Sun S, He J, Zhang Y, Xiao R, Yan M, Ren Y, Zhu Y, Jin T, Xia Y. Genetic

- 1 polymorphisms in the ALDH2 gene and the risk of ischemic stroke in a Chinese han
2 population. *Oncotarget*. 2017;8:101936-101943. doi: 10.18632/oncotarget.21803.
- 3 48. Yang Q, Yin RX, Cao XL, Wu DF, Chen WX, Zhou YJ. Association of two
4 polymorphisms in the FADS1/FADS2 gene cluster and the risk of coronary artery disease
5 and ischemic stroke. *Int J Clin Exp Pathol*. 2015;8:7318-7331.
- 6 49. Gaęało I, Rusiecka I, Kocić I. Tyrosine Kinase Inhibitor as a new Therapy for Ischemic
7 Stroke and other Neurologic Diseases: is there any Hope for a Better Outcome? *Curr*
8 *Neuropharmacol*. 2015;13:836-844. doi: 10.2174/1570159x13666150518235504.
- 9 50. Yokota N, Uchijima M, Nishizawa S, Namba H, Koide Y. Identification of differentially
10 expressed genes in rat hippocampus after transient global cerebral ischemia using
11 subtractive cDNA cloning based on polymerase chain reaction. *Stroke*. 2001;32:168-174.
12 doi: 10.1161/01.str.32.1.168.
- 13 51. Peng H, Fan Y, Li J, Zheng X, Zhong C, Zhu Z, He Y, Zhang M, Zhang Y. DNA
14 Methylation of the Natriuretic Peptide System Genes and Ischemic Stroke: Gene-Based
15 and Gene Set Analyses. *Neurol Genet*. 2022;8:e679. doi:
16 10.1212/NXG.0000000000000679.
- 17 52. Zhang H, Mo X, Wang A, Peng H, Guo D, Zhong C, Zhu Z, Xu T, Zhang Y. Association
18 of DNA Methylation in Blood Pressure-Related Genes With Ischemic Stroke Risk and
19 Prognosis. *Front Cardiovasc Med*. 2022;9:796245. doi: 10.3389/fcvm.2022.796245.
- 20 53. Sundström J, Söderholm M, Borné Y, Nilsson J, Persson M, Östling G, Melander O,
21 Orho-Melander M, Engström G. Eosinophil Cationic Protein, Carotid Plaque, and
22 Incidence of Stroke. *Stroke*. 2017;48:2686-2692. doi: 10.1161/STROKEAHA.117.018450.

1 54. Zhong C, Yang J, Xu T, Xu T, Peng Y, Wang A, Wang J, Peng H, Li Q, Ju Z, et al. Serum
2 matrix metalloproteinase-9 levels and prognosis of acute ischemic stroke. *Neurology*.
3 2017;89:805-812. doi: 10.1212/WNL.0000000000004257.

4 55. Ponsaerts L, Alders L, Schepers M, de Oliveira RMW, Prickaerts J, Vanmierlo T,
5 Bronckaers A. Neuroinflammation in Ischemic Stroke: Inhibition of cAMP-Specific
6 Phosphodiesterases (PDEs) to the Rescue. *Biomedicines*. 2021;9:703. doi:
7 10.3390/biomedicines9070703.

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

1 TABLE 1| Gene name conversion table

Gene description	Gene symbol
ERCC excision repair 2, TFIIH core complex helicase subunit	ERCC2
flap structure-specific endonuclease	FEN1
HD domain containing 3	HDDC3
translocase of outer mitochondrial membrane 40	TOMM40
trafficking protein particle complex subunit 6A	TRAPPC6A
aldehyde dehydrogenase 2 family member	ALDH2
fatty acid desaturase 2	FADS2
FES proto-oncogene, tyrosine kinase	FES
furin, paired basic amino acid cleaving enzyme	FURIN
mannosidase alpha class 2A member 2	MAN2A2
nectin cell adhesion molecule 2	PVRL2
synaptotagmin 7	SYT7
PH domain containing endocytic trafficking adaptor 1	FAM109A

2

3

4

5

6

7

8

9

10

11

12

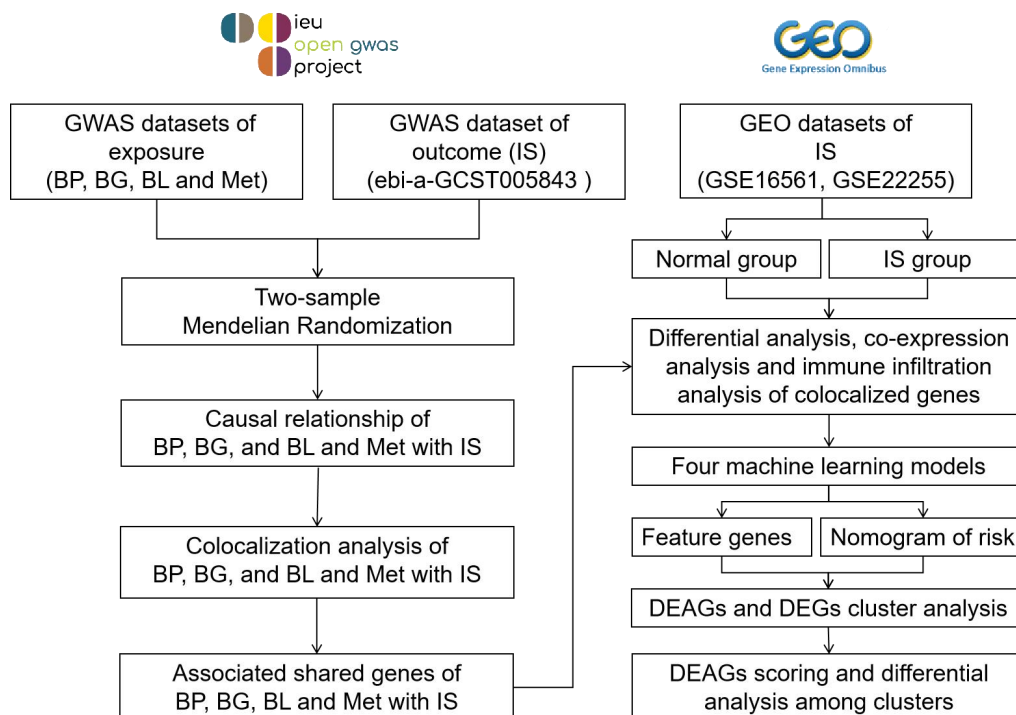
13

14

15

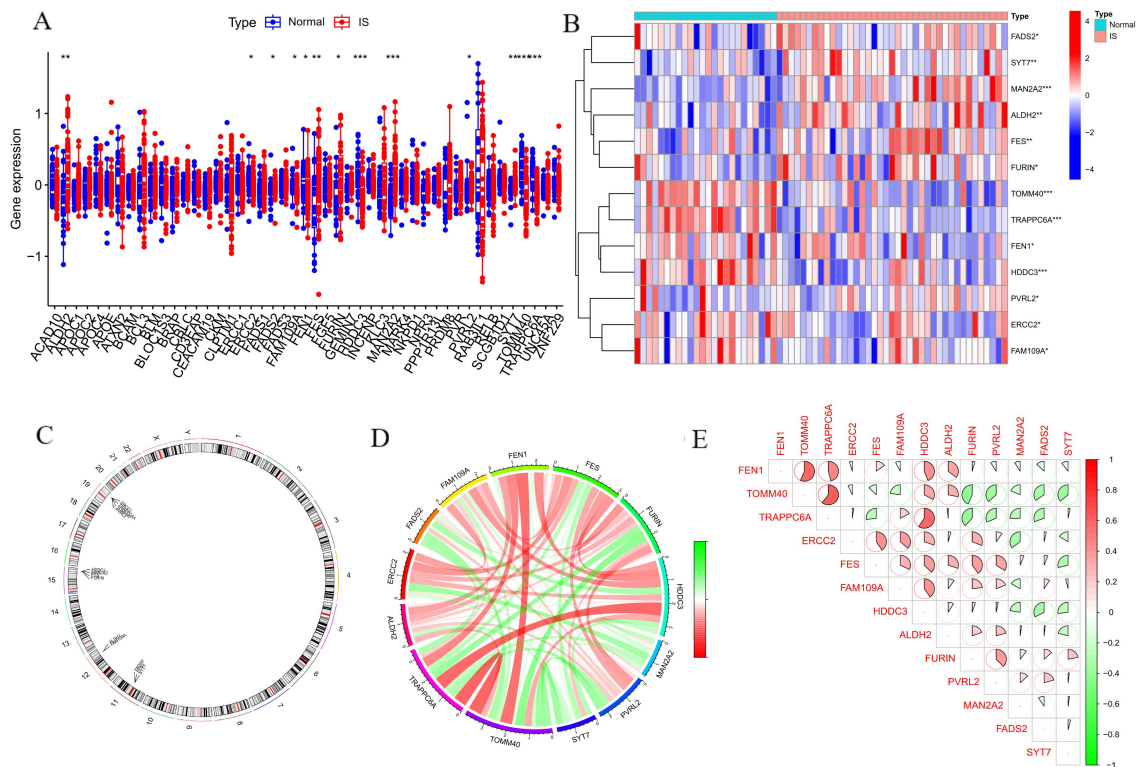
16

1 FIGURE 1 | Flow diagram of this study design



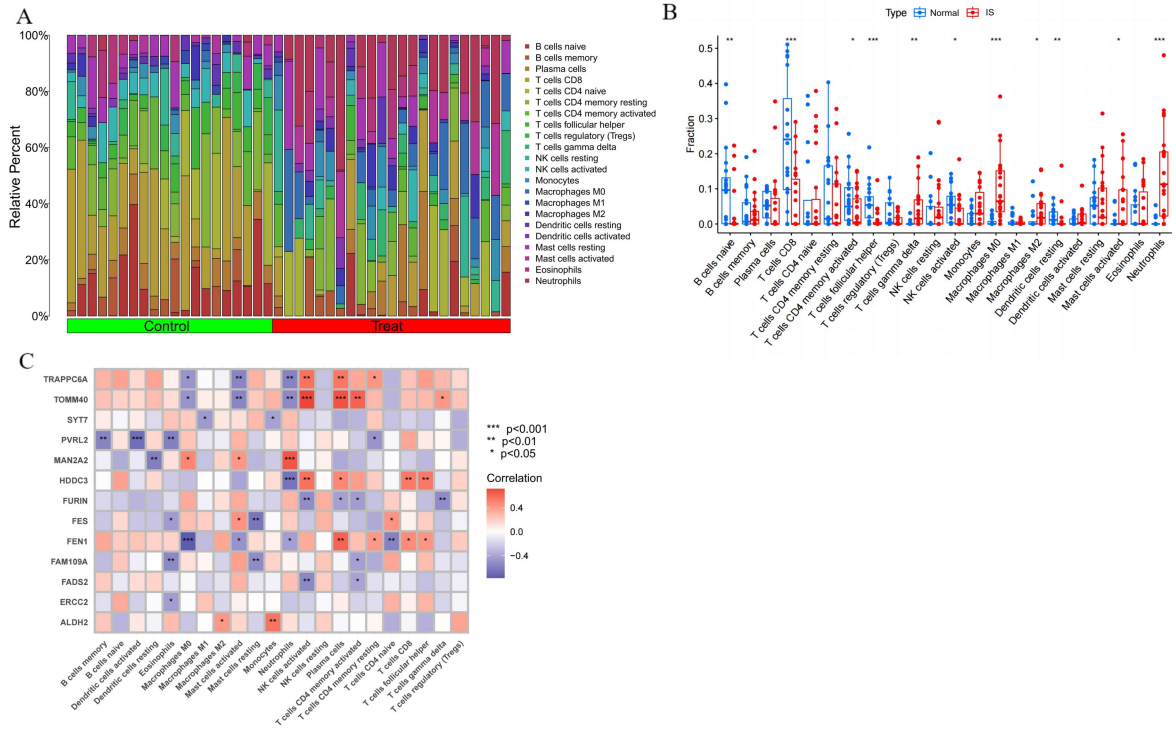
2
3
4
5
6
7
8
9
10
11
12
13
14

1 FIGURE 2 | (A) Box plot of expression difference analysis of associated shared genes
 2 between normal samples and IS samples; (B) Heat map of DEAGs expression in normal and
 3 IS samples; (C) Circle plot of chromosome location of DEAGs; (D) DEAGs correlation
 4 network; (E) Correlation analysis between the two
 5 DEAGs



6
 7
 8
 9
 10
 11
 12
 13

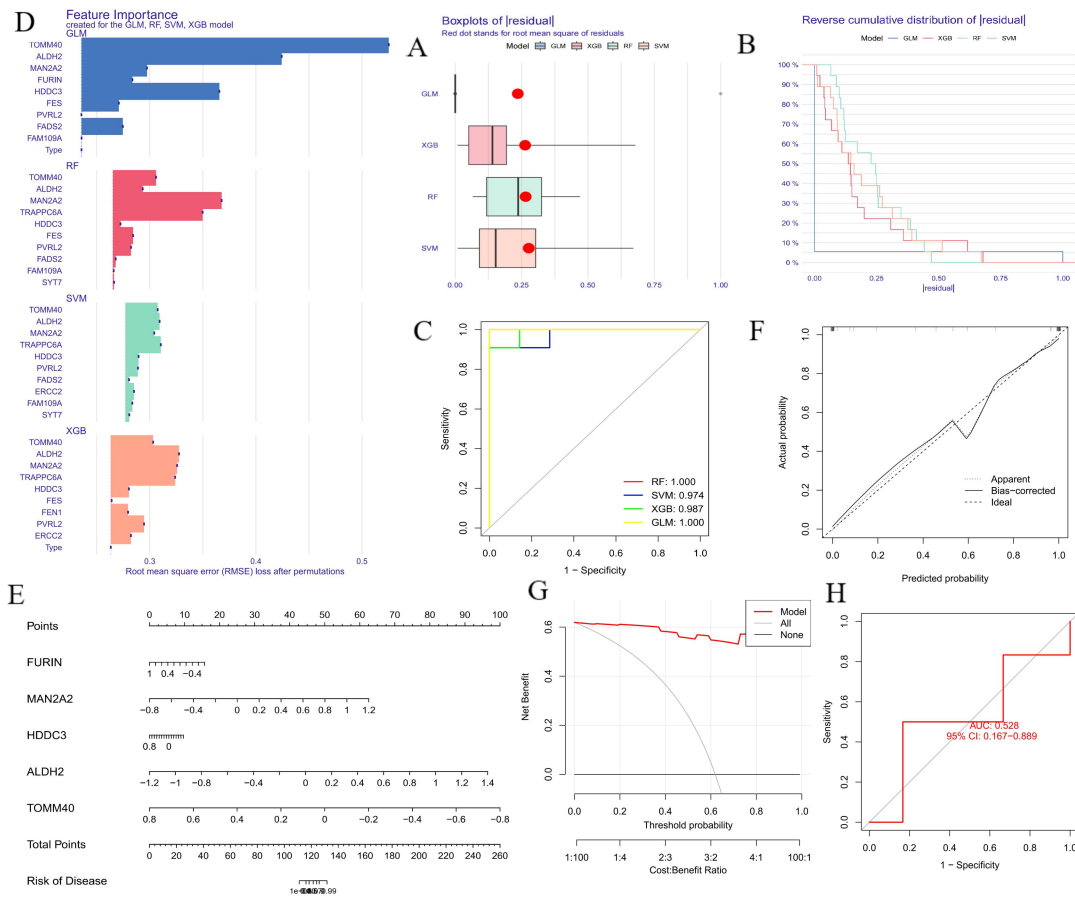
1 FIGURE 3 | (A) Bar plot of relative percentage of each immune cells in samples; (B) Box
 2 plot of immune cell fraction between normal samples and IS samples; (C) Heat map of
 3 correlation analysis between DEAGs and immune cells



4
 5
 6
 7
 8
 9
 10
 11
 12
 13
 14

1 FIGURE 4 | (A) Box plots of residual of the four machine learning models; (B) Reverse
 2 cumulative distribution of residual of the four machine learning models; (C) ROC of the four
 3 machine learning models; (D) Bar plot of feature importance of the four machine learning
 4 models; (E) Nomogram of the feature genes; (F) Calibration curve of feature genes
 5 nomogram; (G) Decision curve of feature genes nomogram; (H) ROC of the test GEO dataset

6



7

8

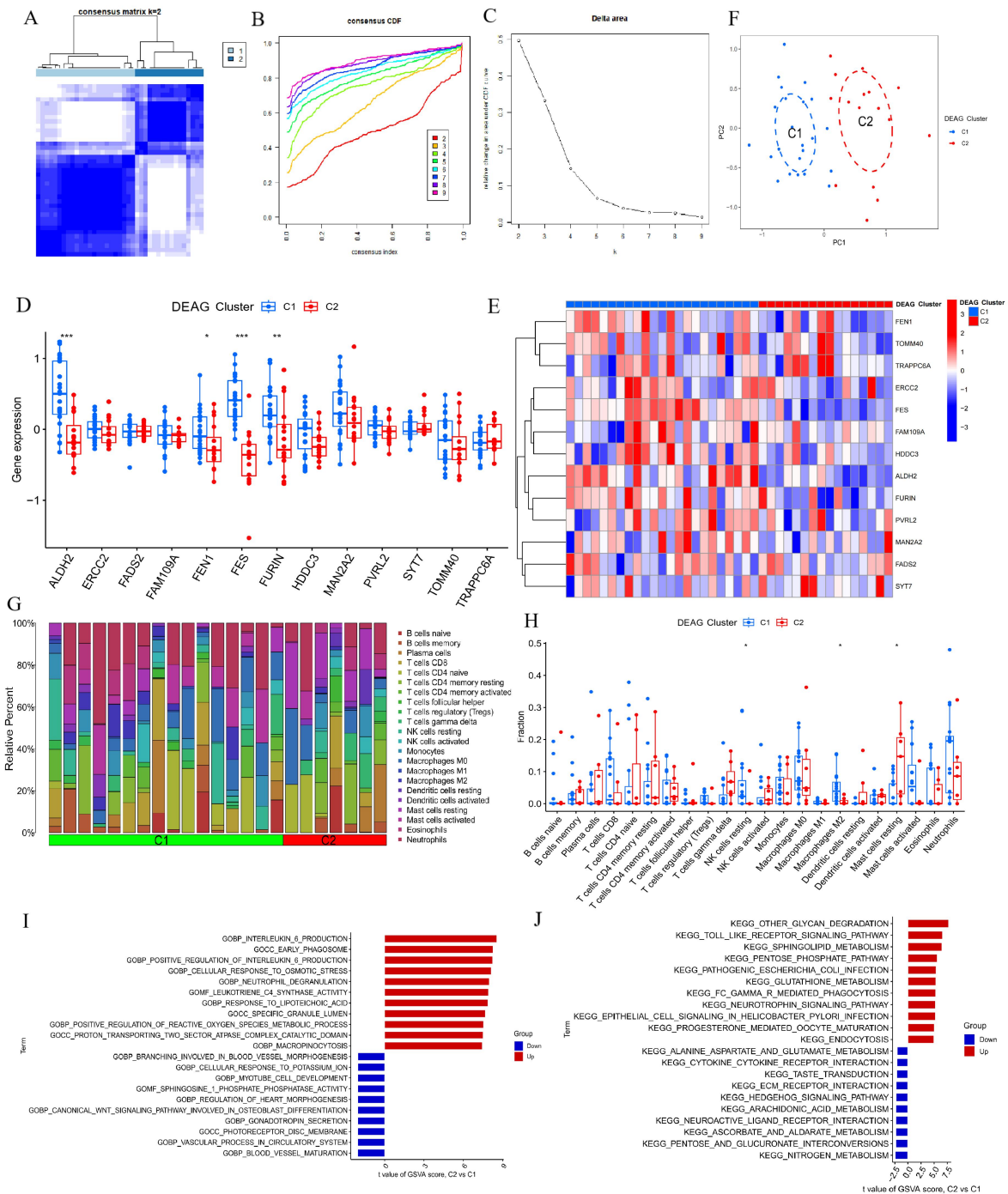
9

10

11

12

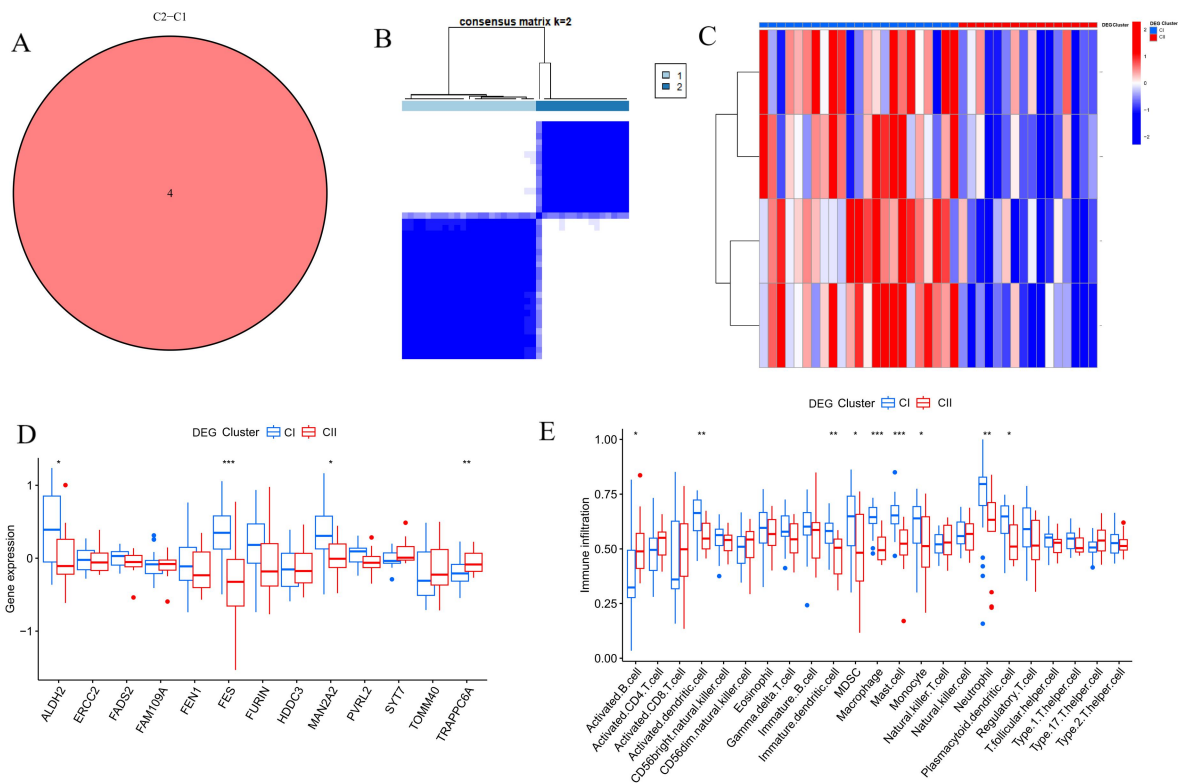
1 FIGURE 5 | (A) Consensus cumulative distribution plot of DEAG clustering for samples; (B)
2 Consensus CDFs of DEAG clustering; (C) Consensus matrix heat map of DEAG clustering
3 for samples; (D) Heat map of DEAGs expression between DEAG clusters; (E) Box plot of
4 expression difference analysis of DEAG clusters; (F) Scatter plot of PCA between DEAG
5 clusters; (G) Box plot of immune cell fraction between DEAG clusters; (H) Bar plot of
6 relative percentage of each immune cells in samples of DEAG clusters; (I) Bar plot of GO
7 terms of GSVA between DEAG clusters; (J) Bar plot of KEGG terms of GSVA between
8 DEAG clusters
9
10
11



1
2
3
4
5
6

1 FIGURE 6 | (A) Venn plot of DEGs; (B) Consensus matrix heat map of DEG clustering for IS
 2 samples; (C) Heat map of DEGs expression between DEG clusters; (D) Box plot of
 3 expression difference analysis of DEG clusters; (E) Box plot of immune cell infiltration
 4 between DEG clusters

5



6

7

8

9

10

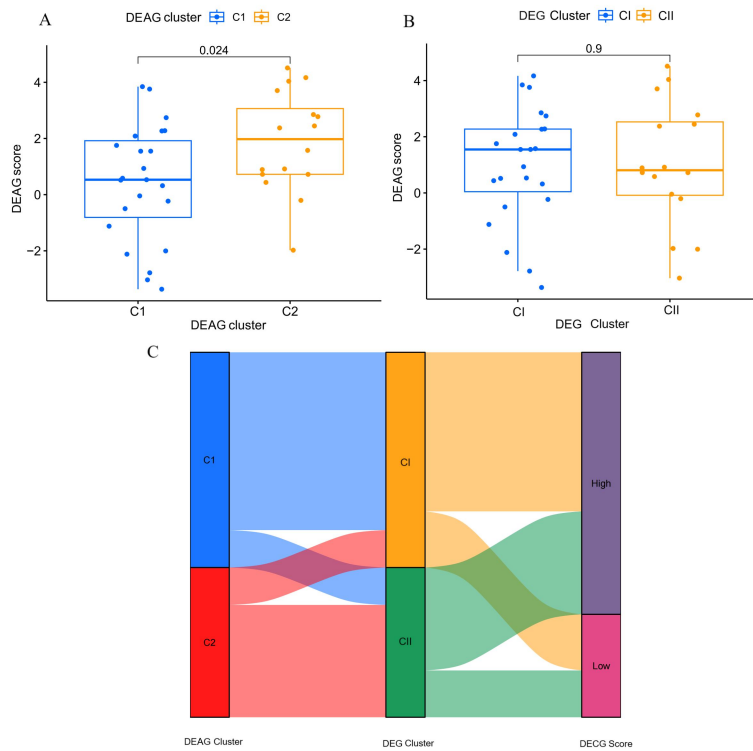
11

12

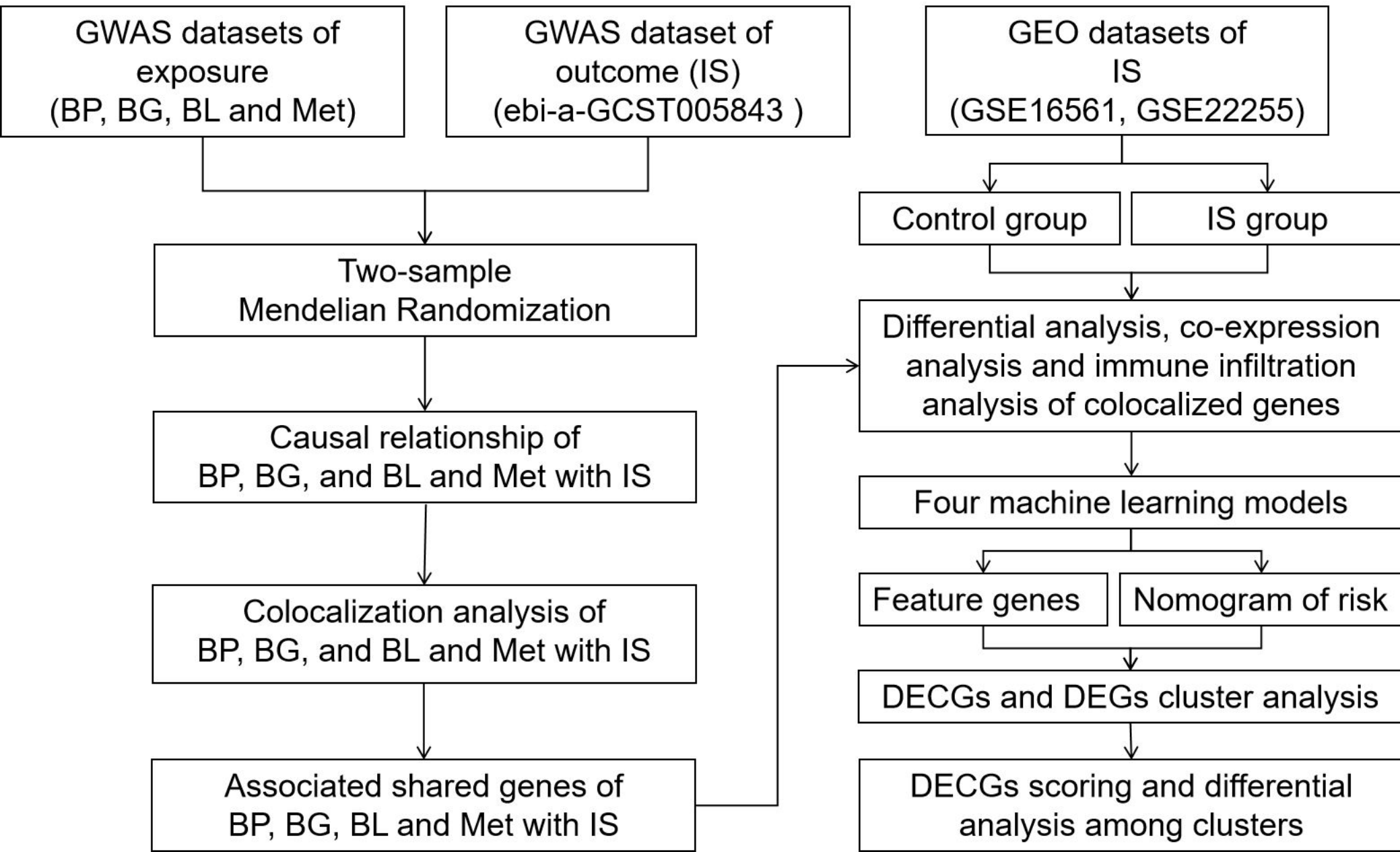
13

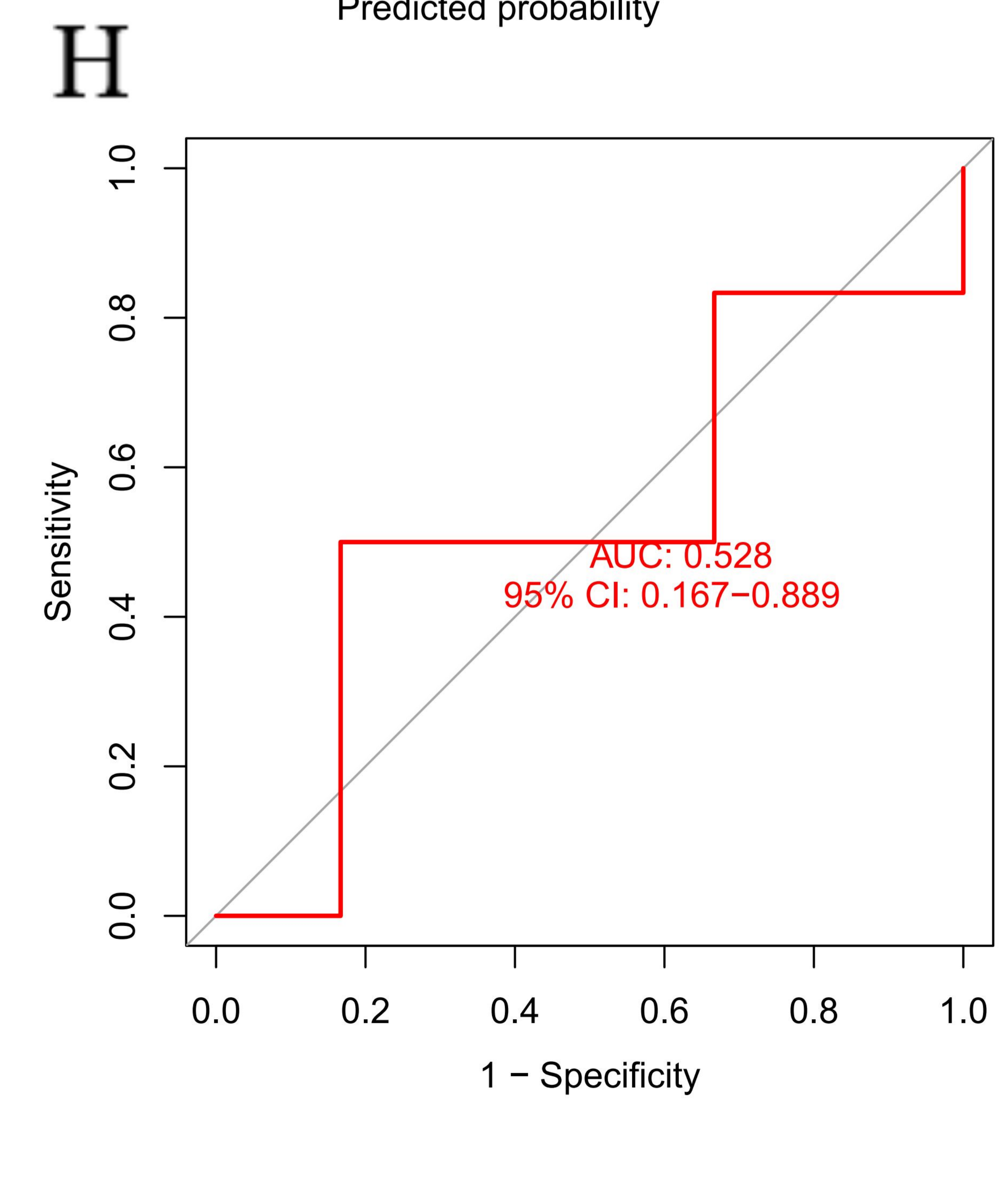
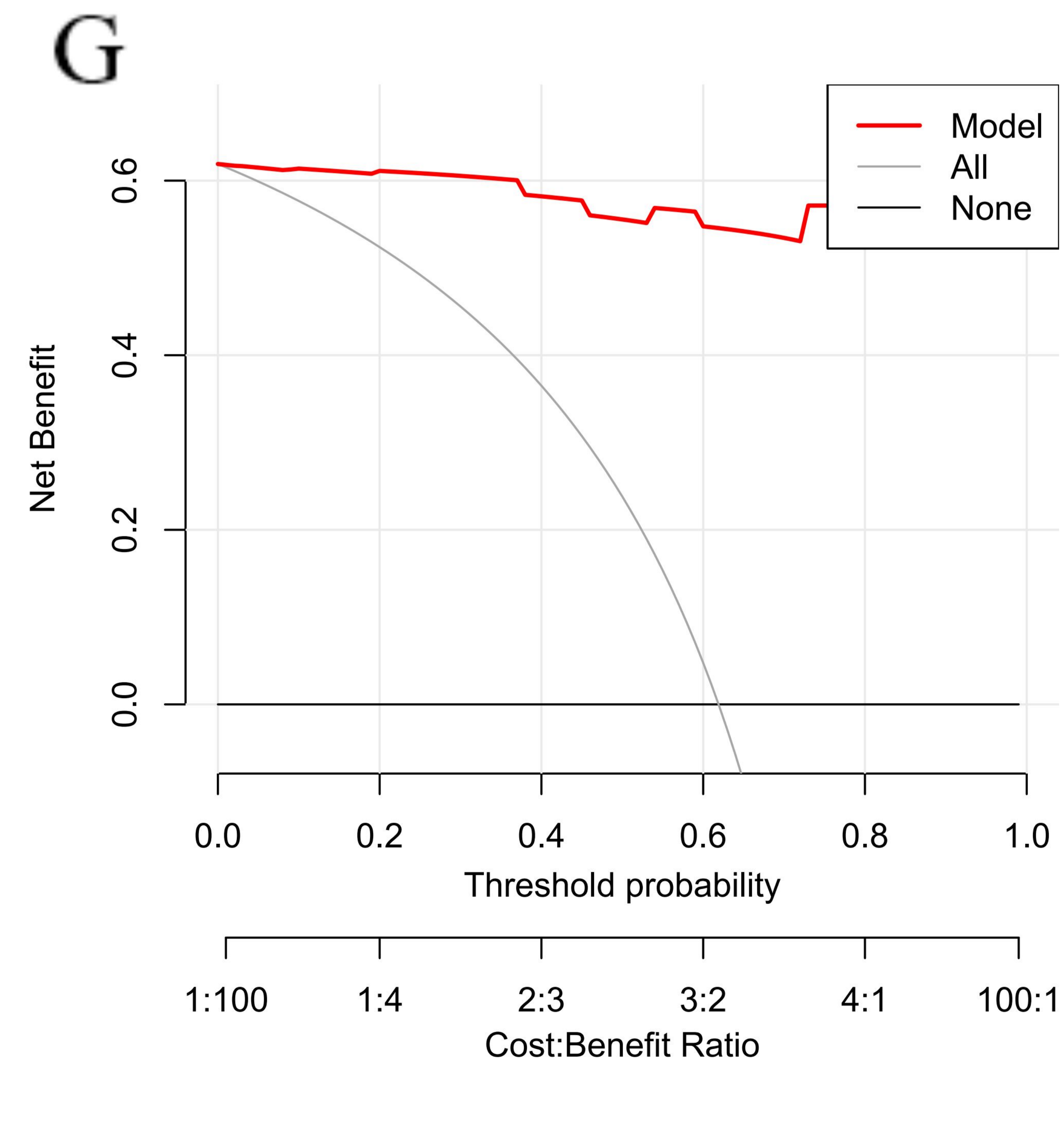
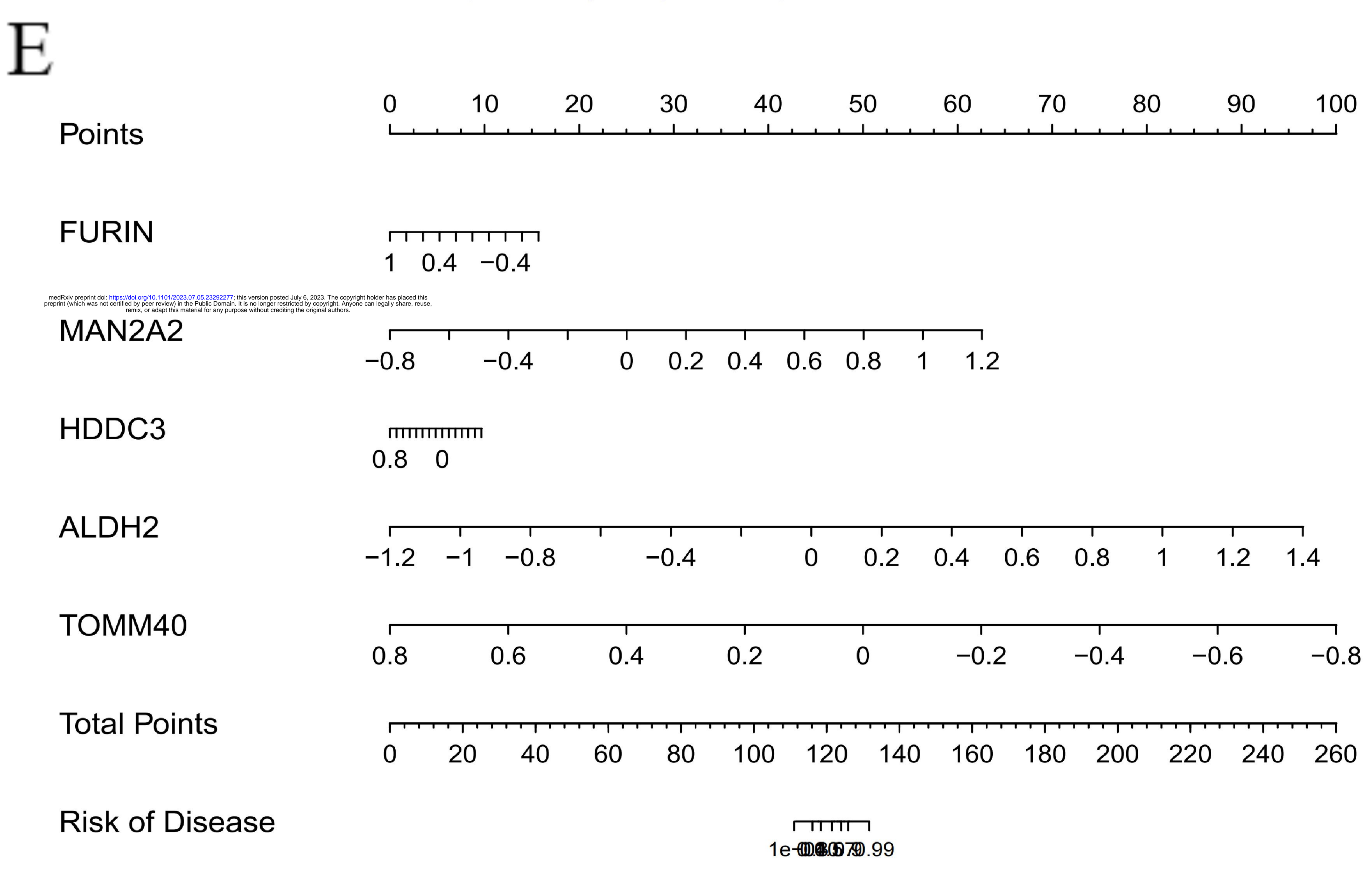
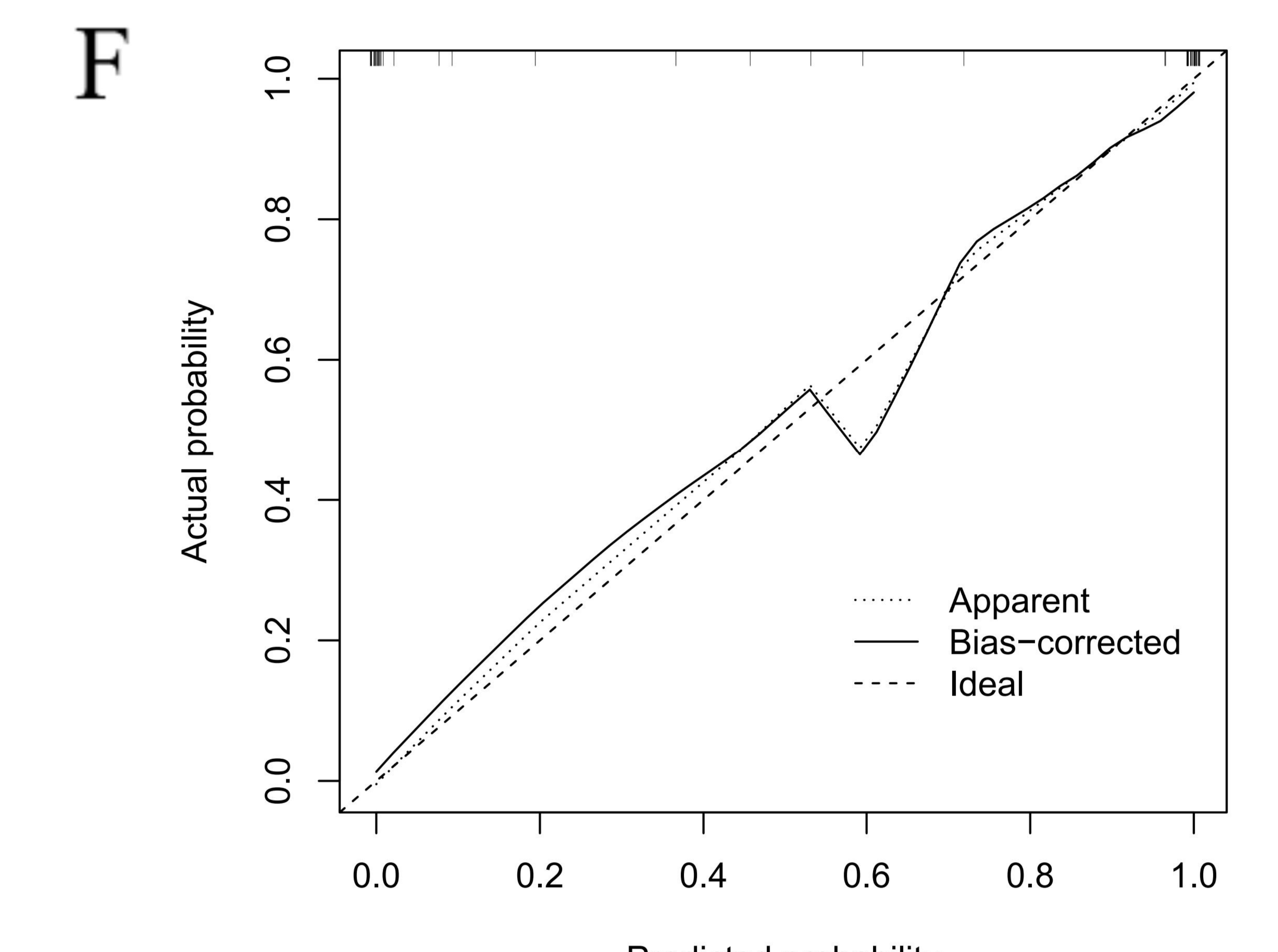
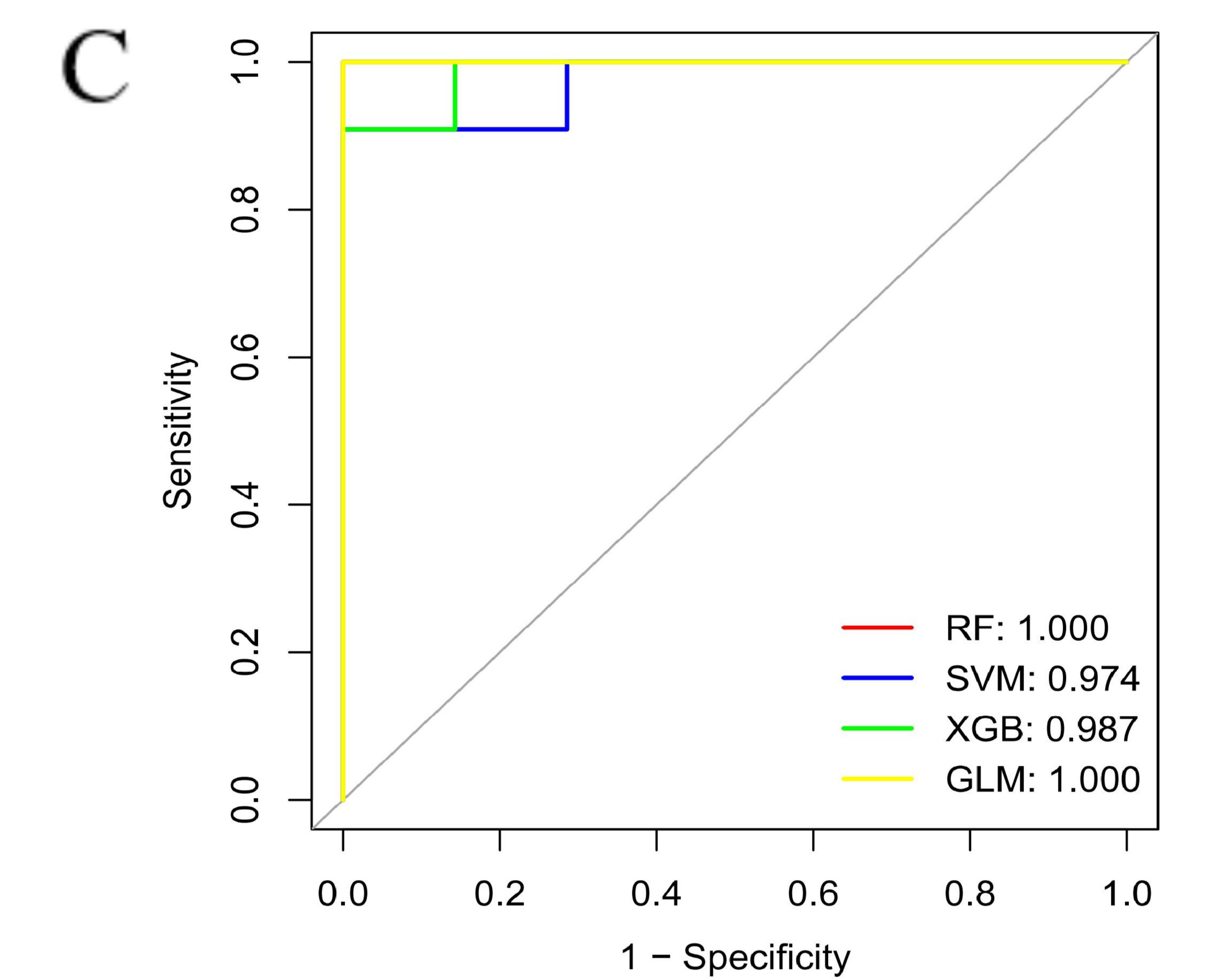
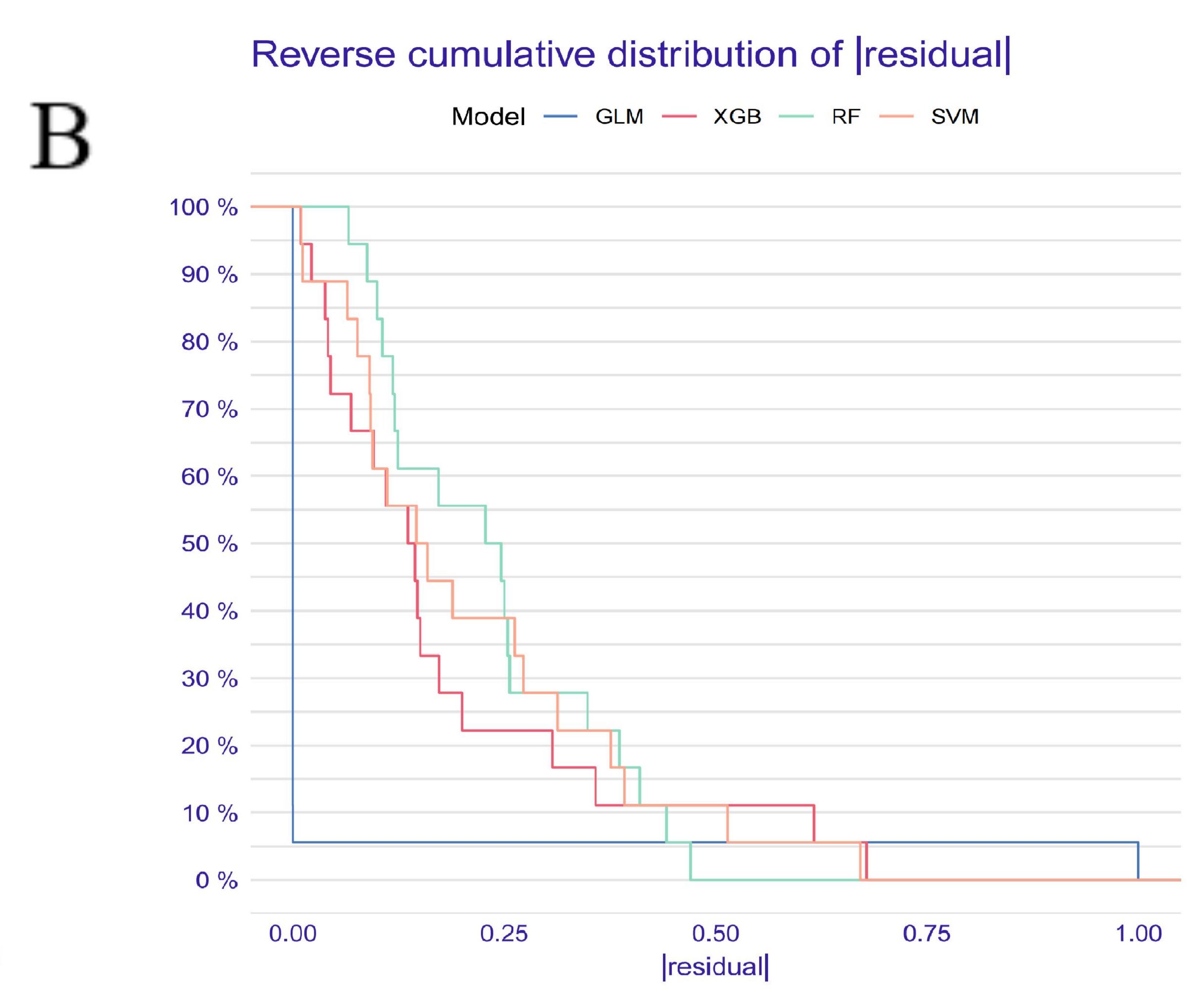
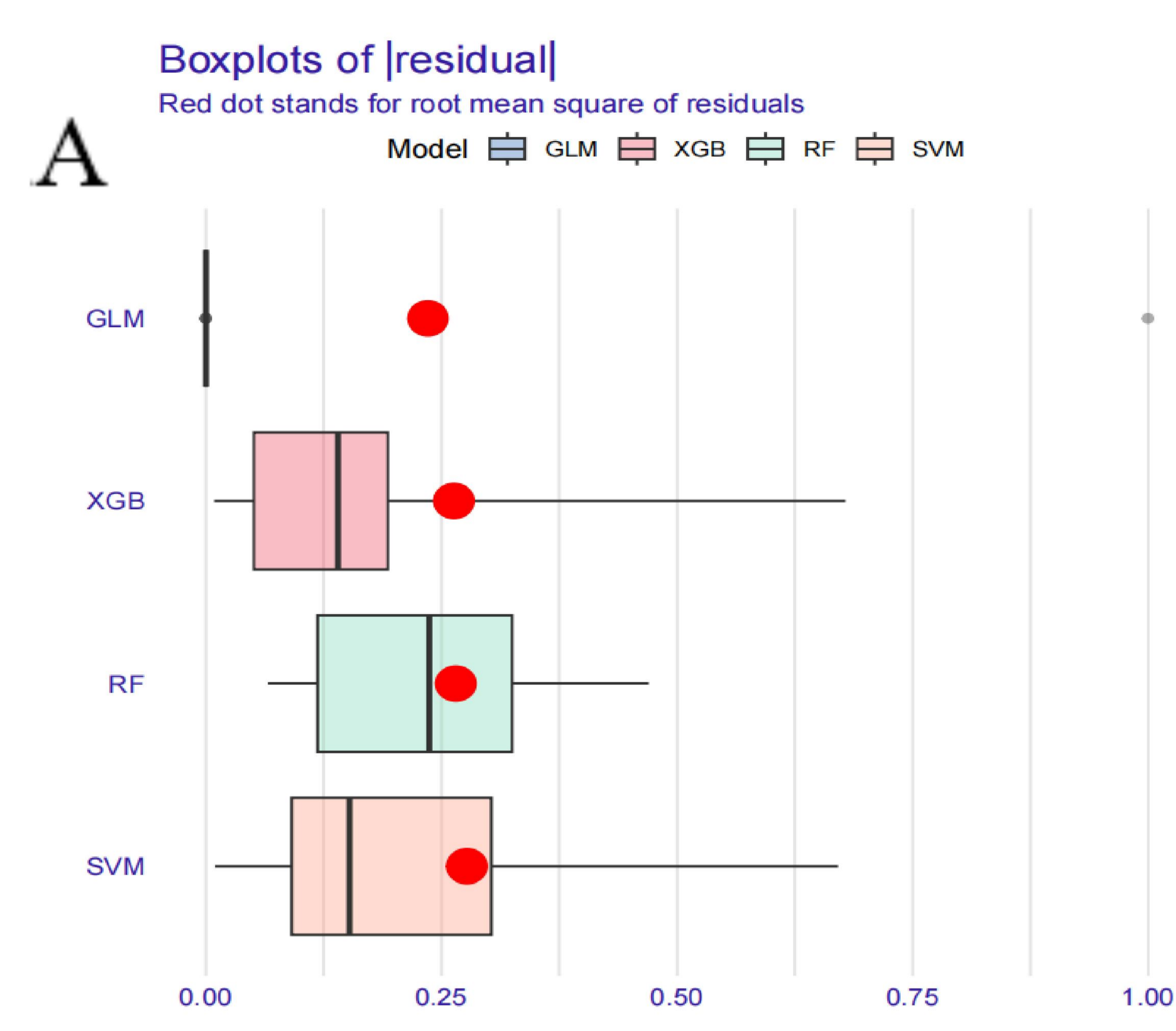
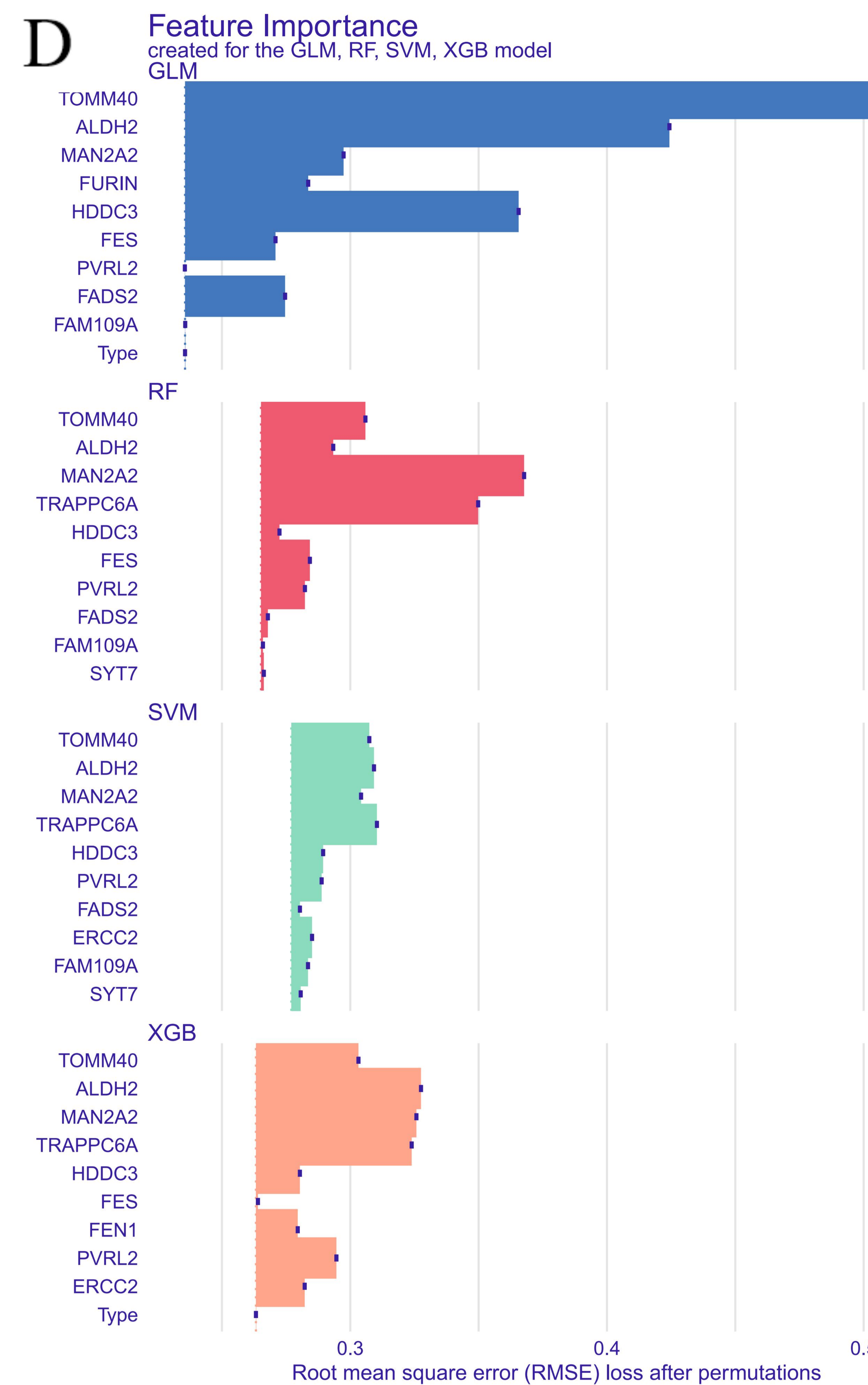
- 1 FIGURE 7 | (A) Box plot of different expression analysis of DEAG score between DEAG
- 2 clusters; (B) Box plot of different expression analysis of DEAG score between DEG clusters;
- 3 (C) Alluvial plot of the correspondence of the different clustered samples

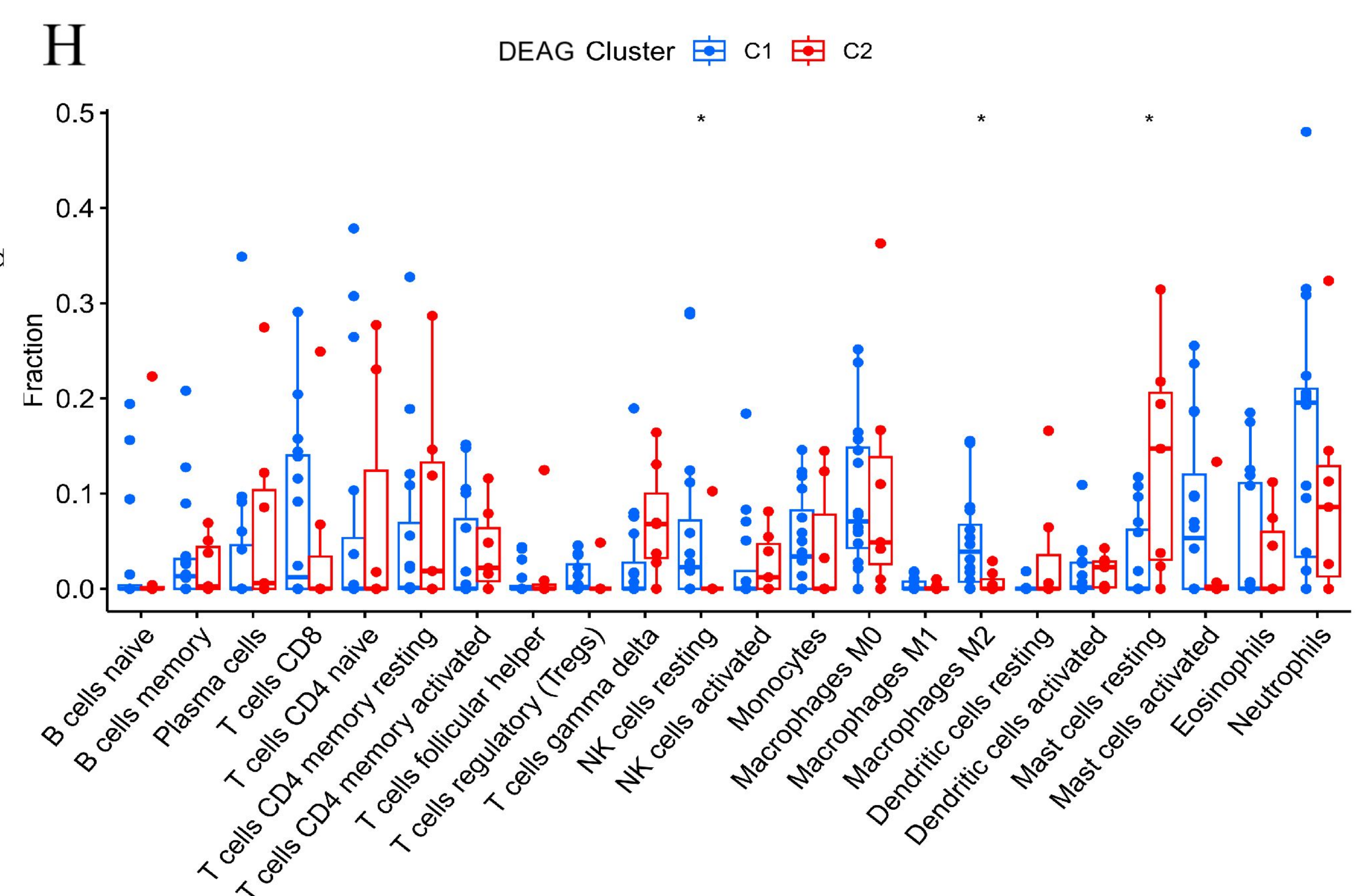
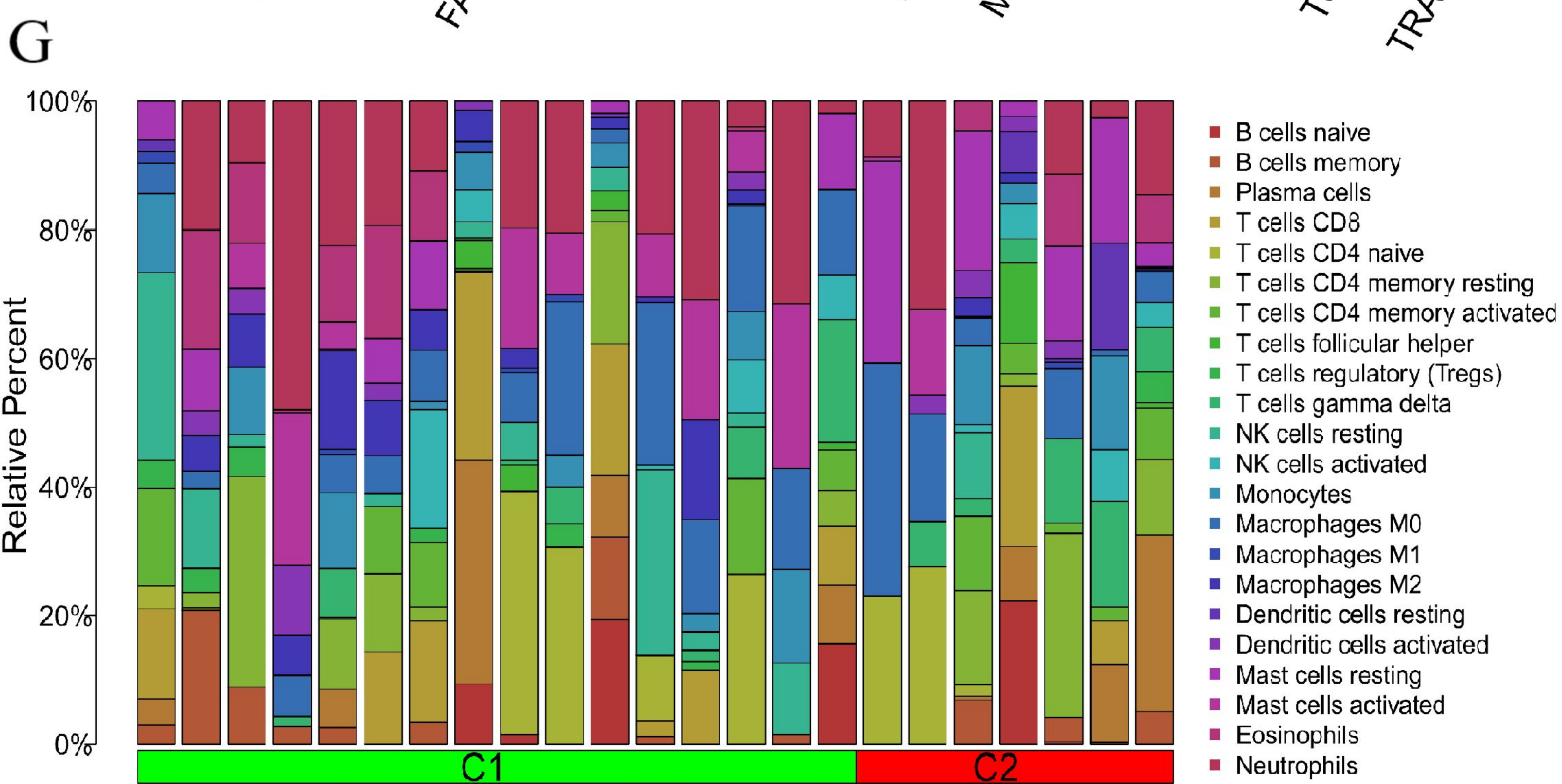
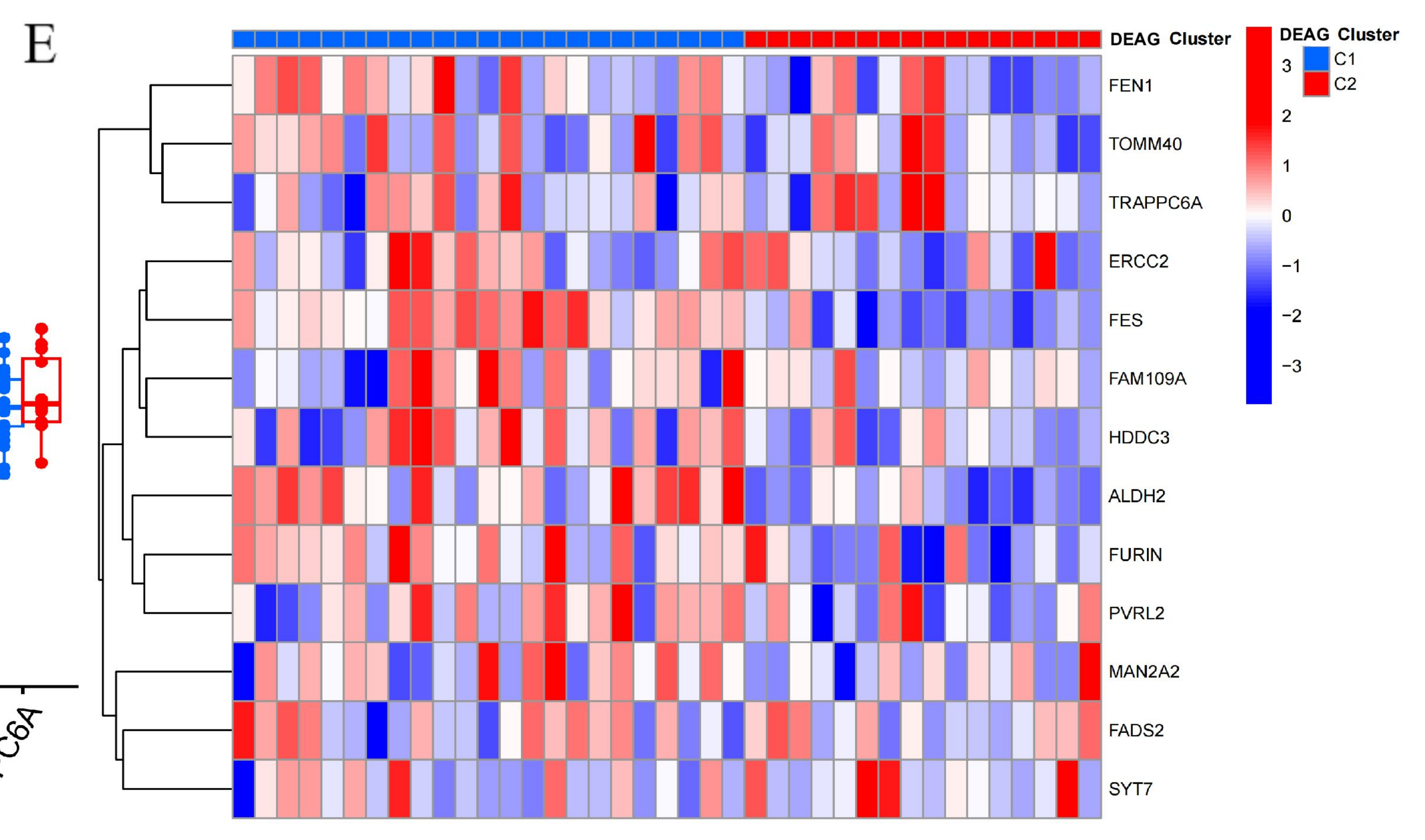
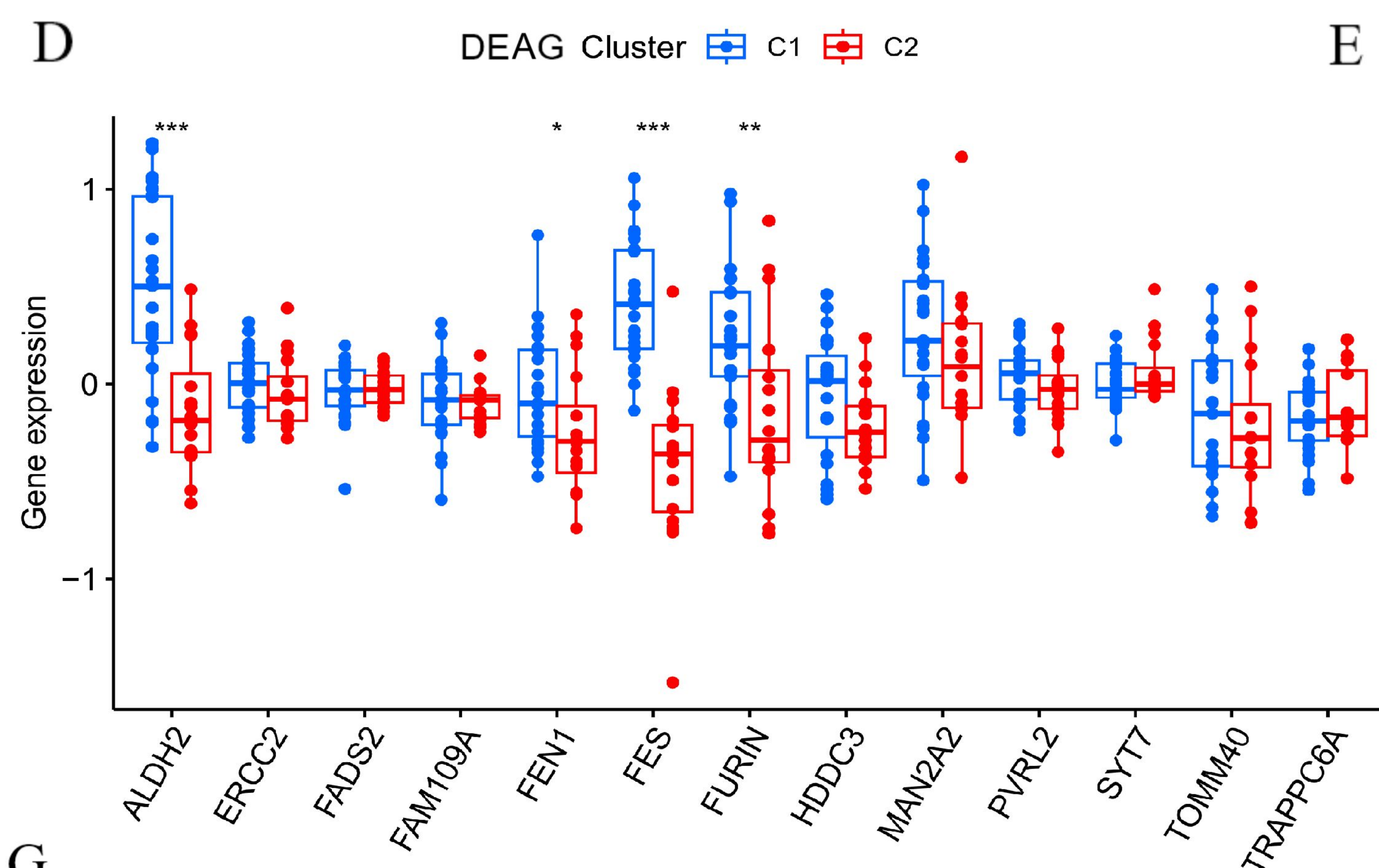
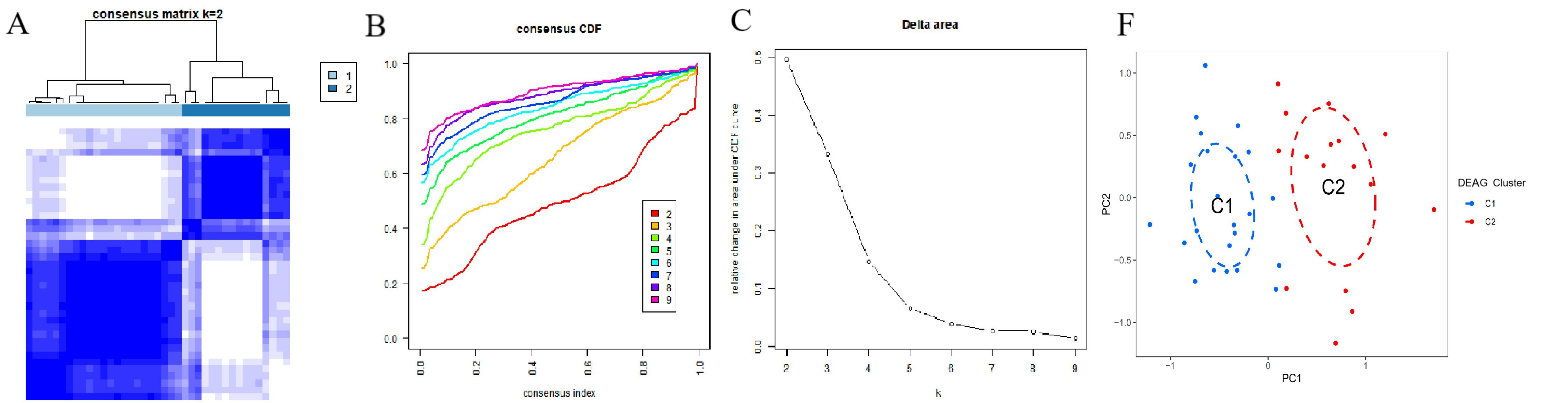
4



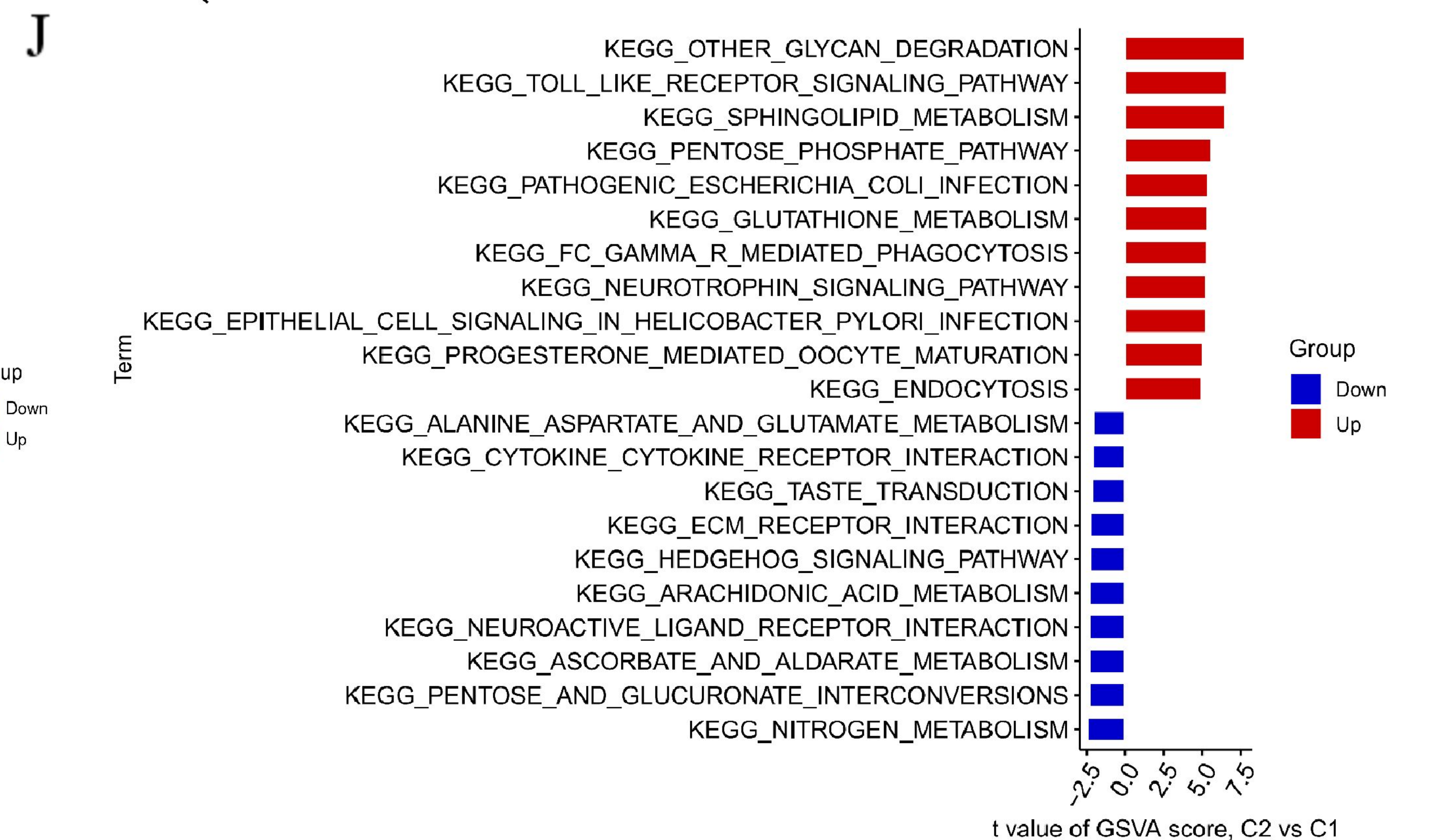
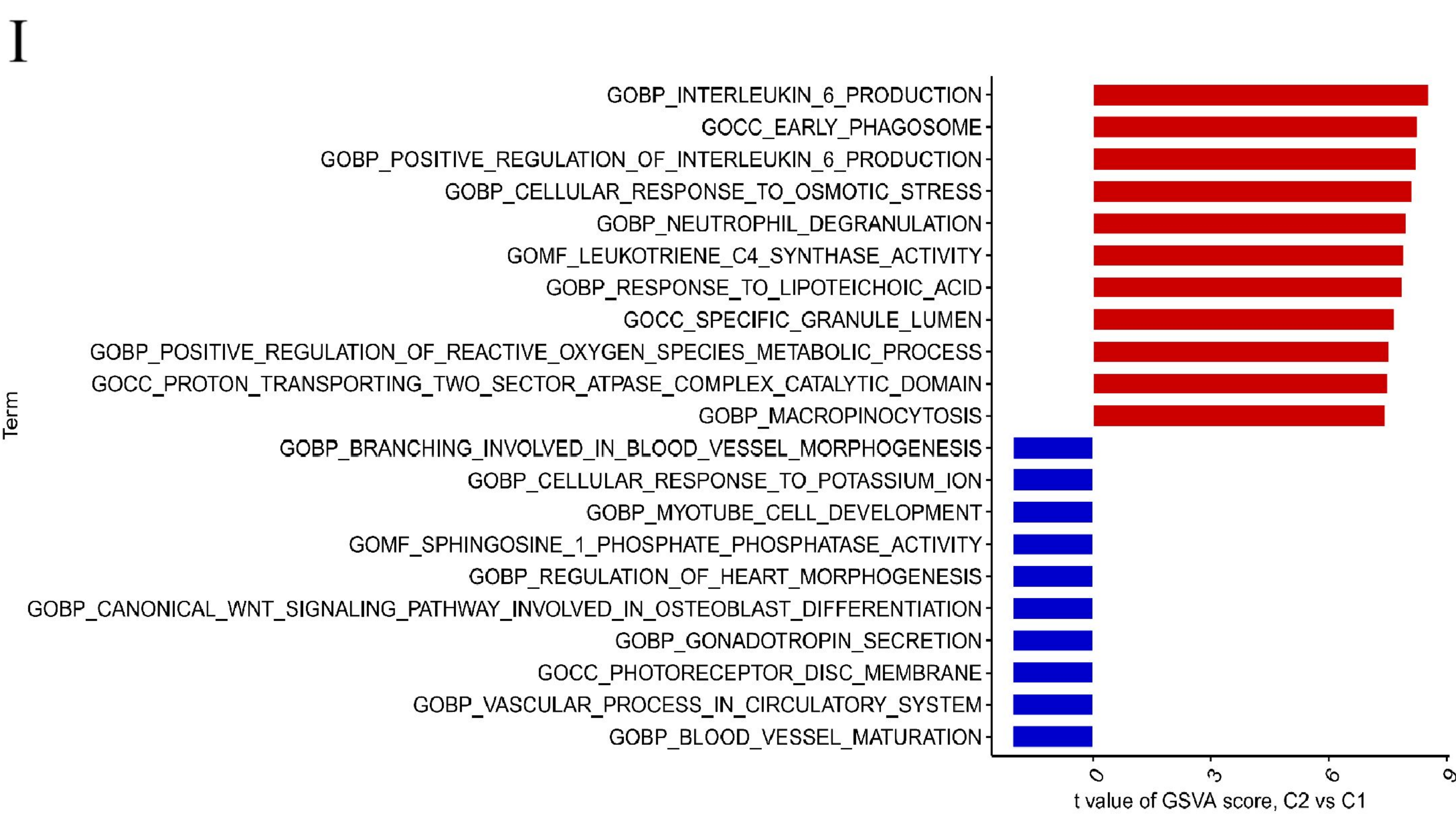
5



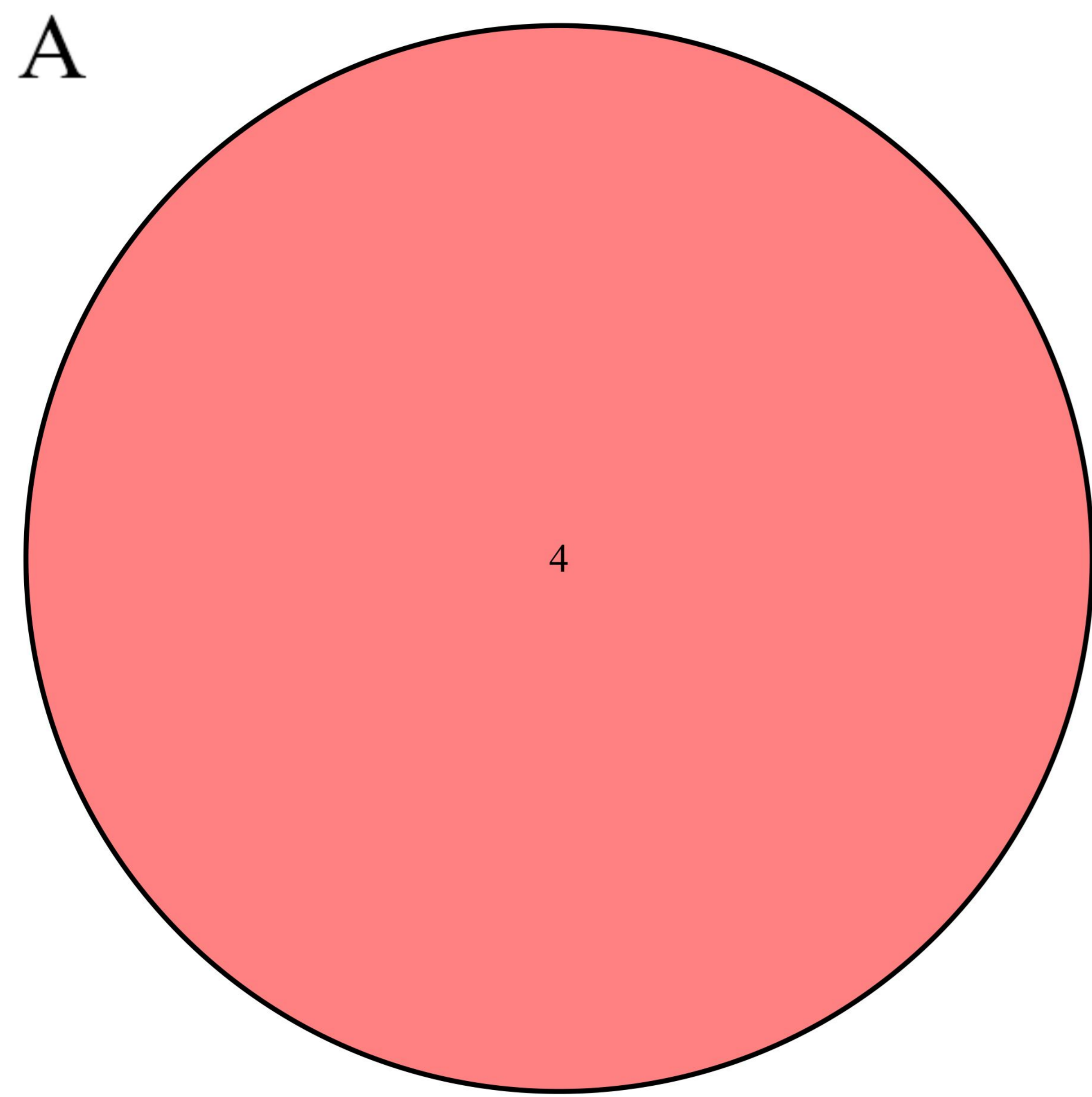
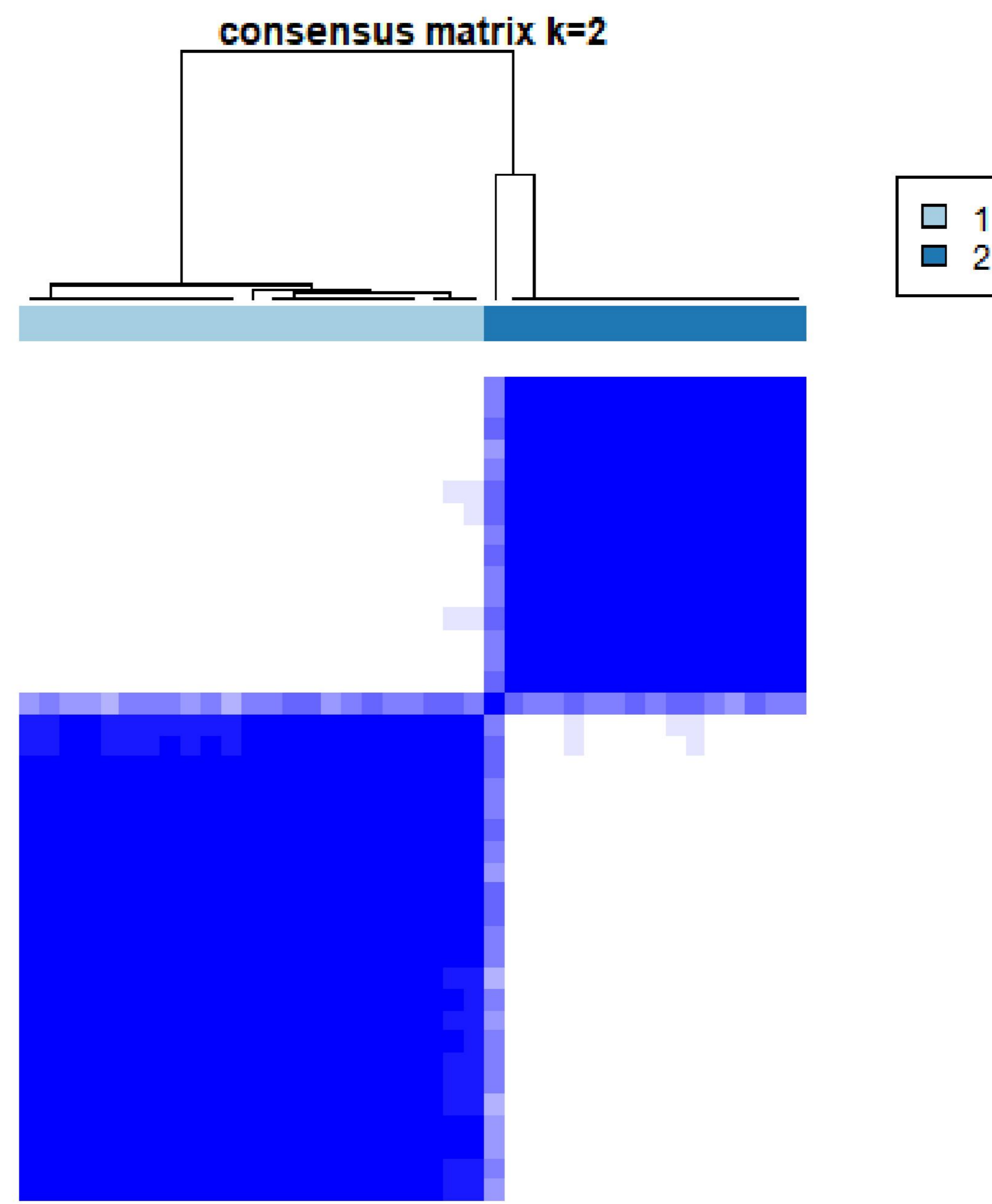
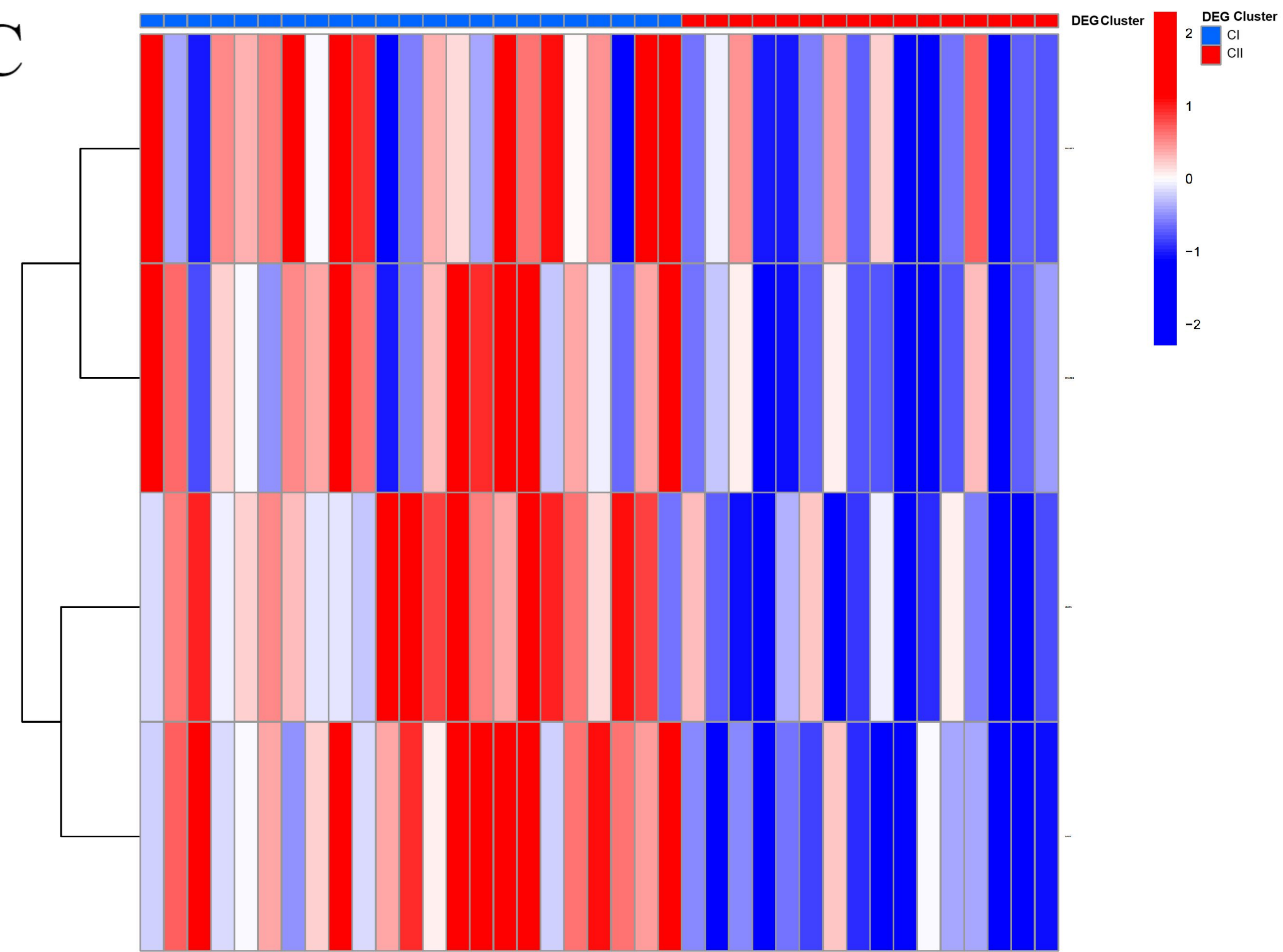
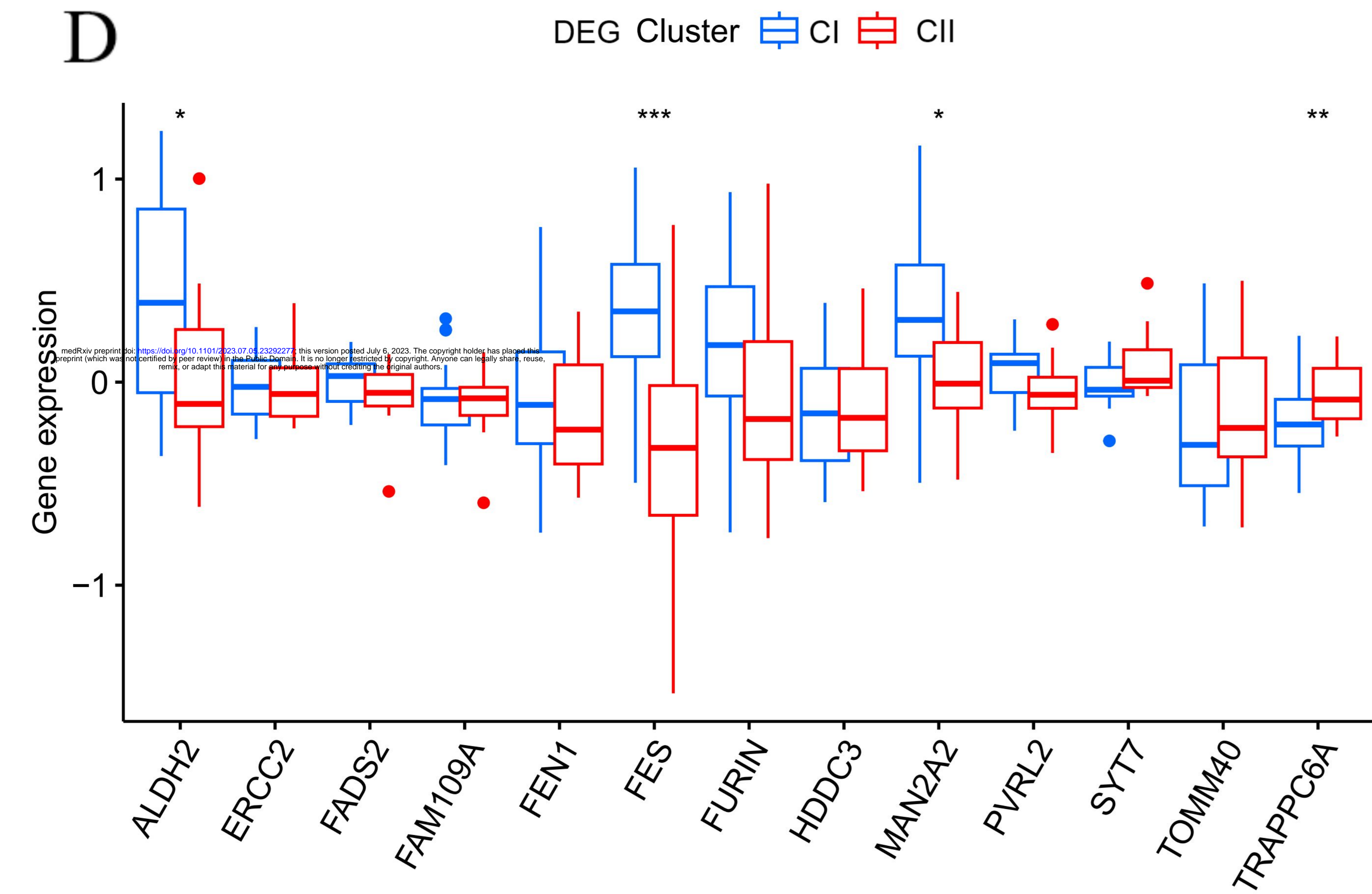
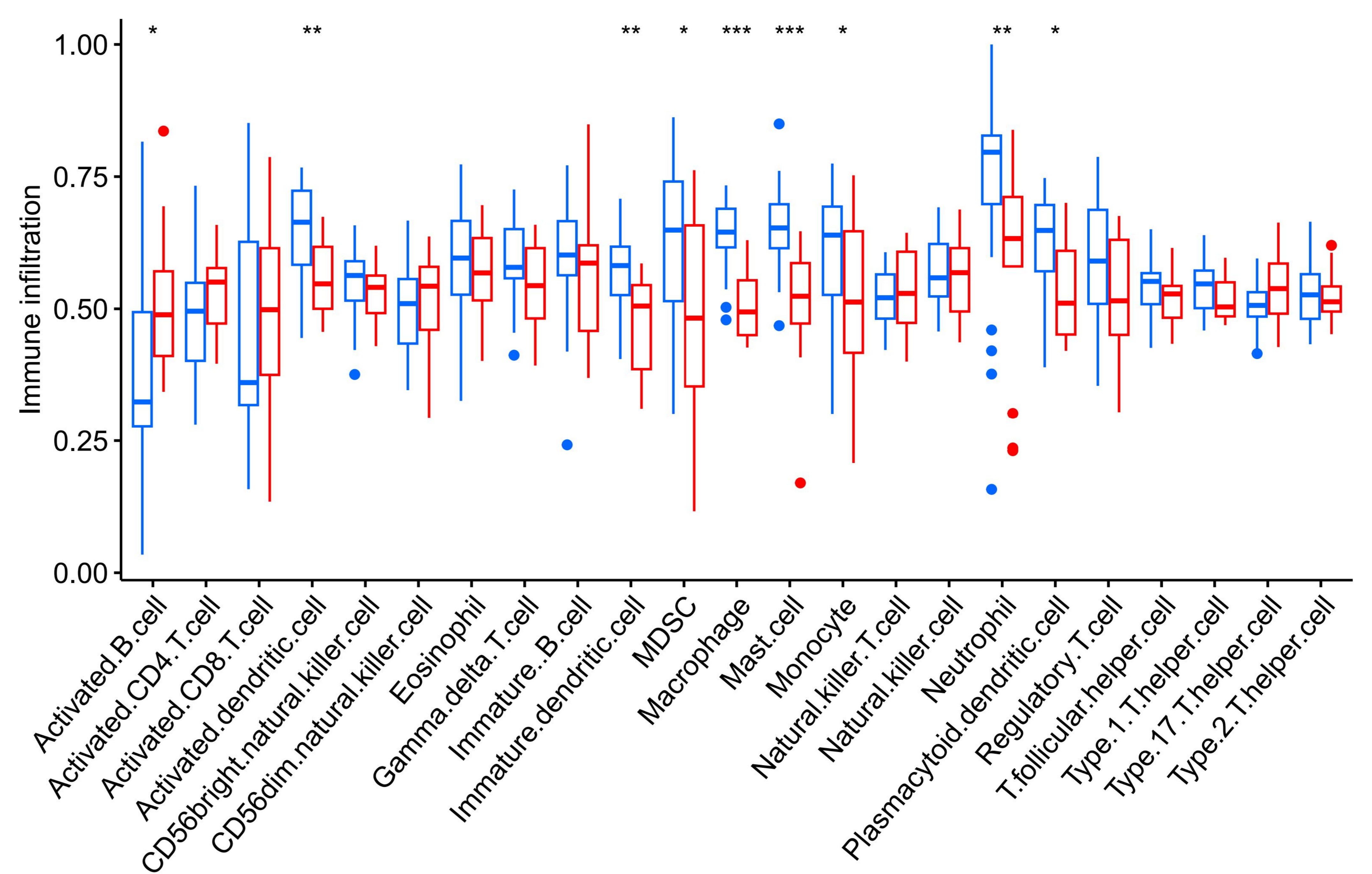




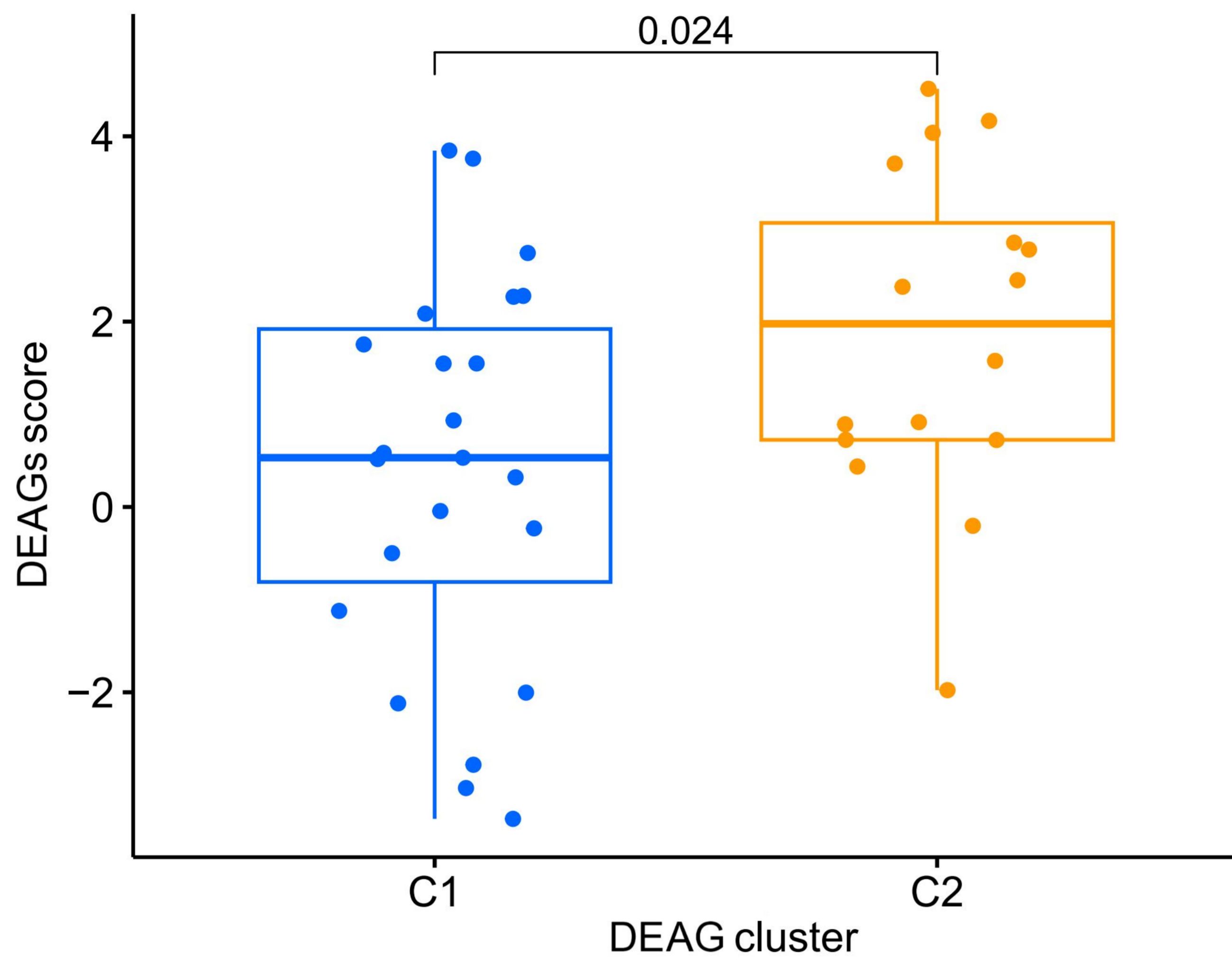
medRxiv preprint doi: <https://doi.org/10.1101/2023.07.05.23292277>; this version posted July 6, 2023. The copyright holder for this preprint (which was not certified by peer review) in the Public Domain. It is no longer restricted by copyright. Anyone can legally share, reuse, remix, or adapt this material for any purpose without crediting the original author.



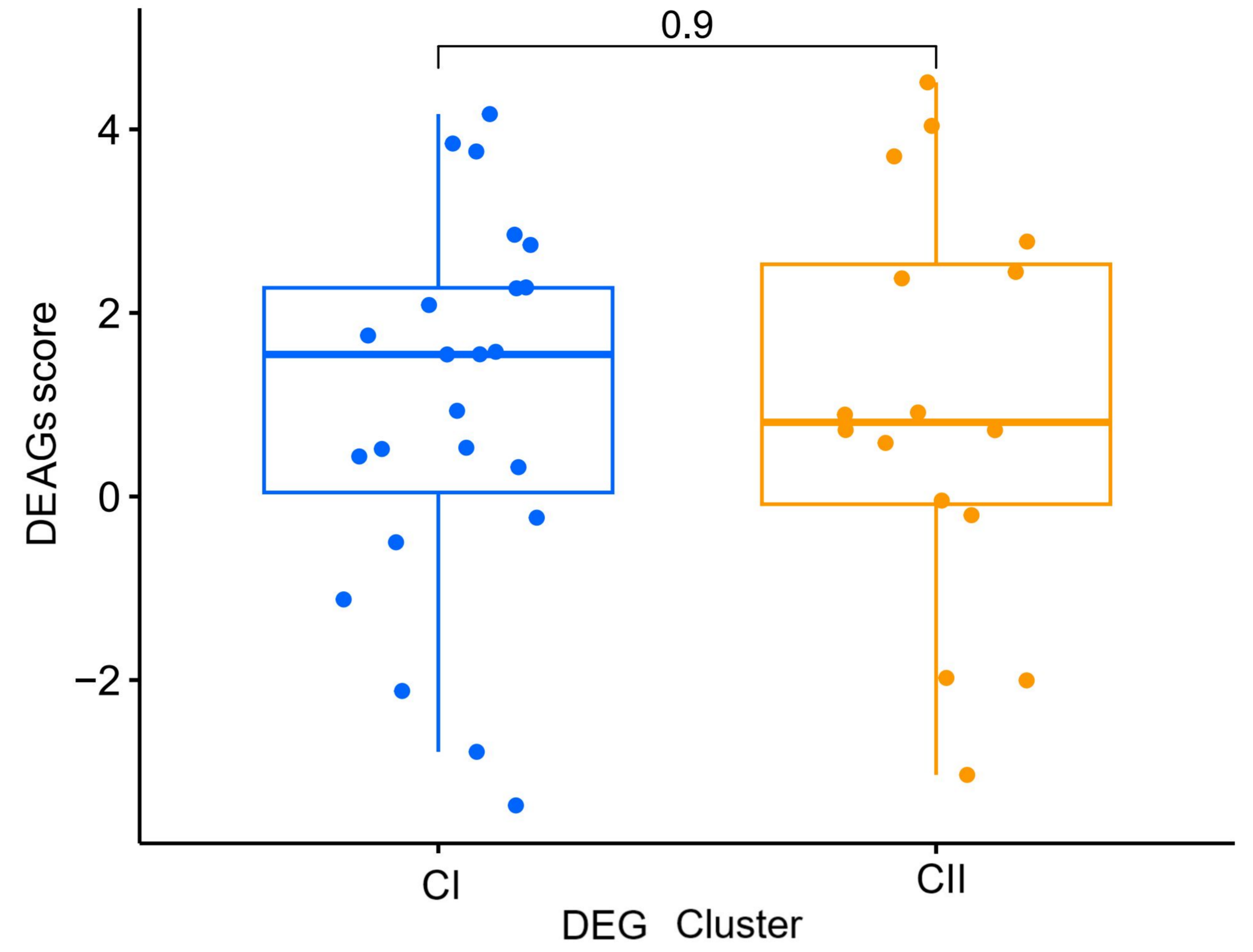
C2-C1

**B****C****D****E**

A

DEAG cluster ■ C1 ■ C2

B

DEG Cluster ■ CI ■ CII

C

medRxiv preprint doi: <https://doi.org/10.1101/2023.07.05.23292277>; this version posted July 6, 2023. The copyright holder for this preprint (which was not certified by peer review) in the Public Domain. It is no longer restricted by copyright. Anyone can legally share, reuse, remix, or adapt this material for any purpose without crediting the original authors.

