## RESEARCH

# Identification of DNA methylation-regulated genes as potential biomarkers for coronary heart disease via machine learning in the Framingham Heart Study

Xiaokang Zhang[1], Chen Wang[1], Dingdong He[1,2], Yating Cheng[1], Li Yu[1], Daoxi Qi[1], Boyu Li[1] and Fang Zheng[1*]

## Abstract

**Background:** DNA methylation-regulated genes have been demonstrated as the crucial participants in the occurrence of coronary heart disease (CHD). The machine learning based on DNA methylation-regulated genes has tremendous potential for mining non-invasive predictive biomarkers and exploring underlying new mechanisms of CHD.

**Results:** First, the 2085 age-gender-matched individuals in Framingham Heart Study (FHS) were randomly divided into training set and validation set. We then integrated methylome and transcriptome data of peripheral blood leukocytes (PBLs) from the training set to probe into the methylation and expression patterns of CHD-related genes. A total of five hub DNA methylation-regulated genes were identified in CHD through dimensionality reduction, including ATG7, BACH2, CDKN1B, DHCR24 and MPO. Subsequently, methylation and expression features of the hub DNA methylation-regulated genes were used to construct machine learning models for CHD prediction by LightGBM, XGBoost and Random Forest. The optimal model established by LightGBM exhibited favorable predictive capacity, whose AUC, sensitivity, and specificity were 0.834, 0.672, 0.864 in the validation set, respectively. Furthermore, the methylation and expression statuses of the hub genes were verified in monocytes using methylation microarray and transcriptome sequencing. The methylation statuses of ATG7, DHCR24 and MPO and the expression statuses of ATG7, BACH2 and DHCR24 in monocytes of our study population were consistent with those in PBLs from FHS.

**Conclusions:** We identified five DNA methylation-regulated genes based on a predictive model for CHD using machine learning, which may clue the new epigenetic mechanism for CHD.

**Keywords:** Coronary heart disease, DNA methylation, Transcriptome, Monocyte, Machine learning

## Introduction

On the basis of data released by the American Heart Association in 2022, the mortality of coronary heart disease (CHD) has reached 1.2‰ in USA [1]. The high mortality of CHD is partly due to the insufficiency and insensitivity of traditional clinical prediction indicators, such as age, gender, serum lipid and blood pressure. For instance, serum cholesterol level failed to identify 2/3 CHD patients in the Framingham Heart Study (FHS) [2]. Even the most extensively used Framingham 10-year

*Correspondence: zhengfang@whu.edu.cn

[1] Center for Gene Diagnosis and Department of Clinical Laboratory Medicine, Zhongnan Hospital of Wuhan University, Donghu Road 169, Wuhan 430071, China

Full list of author information is available at the end of the article

Zhang *et al. Clinical Epigenetics*    (2022) 14:122

Page 2 of 17

risk scale exhibited poor accuracy in CHD prediction [3]. Mining more sensitive non-invasive biomarkers will be of great benefit to the population screening and prediction of CHD, as well as epigenetic pathogenesis discovery.

DNA methylation has been widely demonstrated to be involved in the pathogenesis of CHD [4, 5]. Among the mechanisms by which DNA methylation affects the progression of CHD, the regulation of gene expression by promoter methylation plays a major role. For instance, the hypermethylation of brain and muscle ARNT-like protein 1 (BMAL1) promoter was proved to suppress BMAL1 transcription and further elevate oxidative stress and inflammatory response in human aortic endothelial cells [6]. Like BMAL1, genes regulated by their methylation status are defined as DNA methylation-regulated genes [7]. Identification of CHD-associated DNA methylation-regulated genes may shed light on the pathogenesis of CHD and improve the accuracy of CHD prediction.

In the era where the barrier to profile and access high-throughput data with large sample size was quite low, machine learning is becoming more widely applied in both basic medical research and clinical practice. The essence of machine learning is to generalize features and construct models from existing data through specific algorithms, so that the optimized models can be applied in the prediction of new data [8–10]. Combining appropriate algorithms with prospective CHD biomarkers to construct machine learning models will be advantageous for hub gene identification and prediction of CHD.

In the present study, we integrated methylome and transcriptome data of 2085 individuals from FHS to screen candidate DNA methylation-regulated genes in CHD. After multiple dimensionality reductions, predictive models were created based on five genes' methylation and/or expression statuses using three machine learning algorithms. Finally, the expression and methylation statuses of the five genes were further confirmed in our population using DNA methylation microarray and transcriptome sequencing.

## Results
### Baseline characteristics of the individuals from FHS
There were 2085 individuals in the FHS with methylome data and transcriptome data met our standards. We randomly divided these individuals into training set and validation set. The clinical baseline information of the two subsets is displayed in Table 1. It could be observed that diastolic blood pressure (DBP) and high-density lipoprotein cholesterol (HDL-C) were much lower in CHD patients than in controls. Consistent with previous studies, higher levels of body mass index (BMI), fasting blood glucose (FBG) and triglycerides (TG) were found in CHD patients. Instead of increasing, low-density lipoprotein cholesterol (LDL-C) and total cholesterol (TC) were significantly reduced in CHD patients. And there were no significant differences in smoking and systolic blood pressure (SBP), which were widely reported to predispose to CHD.

**Table 1** Clinical characteristics of the individuals from FHS

| Characteristic | Training controls | Training CHD patients | *p* value | Validation controls | Validation CHD patients | *p* value |
|---|---|---|---|---|---|---|
| Number | 1364 | 199 | – | 455 | 67 | – |
| Male/Female | 554/810 | 95/104 | 0.0644[c] | 186/269 | 32/35 | 0.2917[c] |
| Age (year) | 66 (60, 72) | 66 (61, 72) | 0.1747[b] | 60 (52, 65) | 63 (53, 67) | 0.2557 [b] |
| Smoker / non-smoker | 101 / 1261 | 19 / 180 | 0.3173[c] | 39 / 416 | 10 / 57 | 0.1138[c] |
| BMI (kg/m$^2$) | 27.55 (24.51, 30.84) | 29.38 (25.30, 32.70) | 0.0003[b] | 27.39 (24.82, 30.90) | 28.97 (26.48, 33.15) | 0.0133[b] |
| FBG (mg/dL) | 101.00 (95.00, 110.00) | 106.00 (97.00, 120.00) | < 0.0001[b] | 101.00 (94.00, 110.30) | 109.00 (99.00, 122.00) | 0.0005[b] |
| SBP (mmHg) | 127.00 (117.00, 139.00) | 130.50 (118.00, 141.30) | 0.1229[b] | 127.00 (116.00, 139.00) | 128.50 (118.00, 139.00) | 0.4886[b] |
| DBP (mmHg) | 74.00 (68.00, 80.00) | 71.00 (66.00, 80.00) | 0.0100[b] | 74.05 ± 10.12 | 70.66 ± 11.14 | 0.0118[a] |
| Hypertension treatment (treated/ untreated) | 596 / 700 | 130 / 56 | < 0.0001[c] | 200 / 225 | 43 / 19 | 0.0011[c] |
| HDL-C (mg/dL) | 57.00 (46.00, 69.00) | 48.00 (40.00, 60.00) | < 0.0001[b] | 58.00 (46.00, 70.00) | 52.00 (41.00, 62.00) | 0.0136[b] |
| LDL-C (mg/dL) | 106.00 (86.00, 126.80) | 89.50 (71.00, 112.30) | < 0.0001[b] | 108.00 ± 30.09 | 92.52 ± 31.68 | 0.0001[a] |
| TC (mg/dL) | 188.00 (164.00, 214.00) | 169.00 (142.00, 195.5) | < 0.0001[b] | 189.20 ± 34.78 | 171.10 ± 38.22 | 0.0002[a] |
| TG (mg/dL) | 101.00 (73.00, 139.00) | 125.00 (81.00, 161.30) | < 0.0001[b] | 99.00 (71.00, 133.50) | 120.50 (82.00, 150.50) | 0.0053[b] |

Data were showed as mean ± standard deviation or median (interquartile range)

*BMI* Body mass index, *FBG* Fasting blood glucose, *SBP* Systolic blood pressure, *DBP* Diastolic blood pressure, *HDL-C* High-density lipoprotein cholesterol, *LDL-C* Low-density lipoprotein cholesterol, *TC* Total cholesterol, *TG* Triglycerides

[a] Student's *t* test. [b] Mann–Whitney *U* test. [c] Chi-square test

Zhang *et al. Clinical Epigenetics*    (2022) 14:122

Page 3 of 17

## DMPs distribution and gene set enrichment pattern in CHD

We identified 40,015 differentially methylated positions (DMPs) between 199 CHD patients and 1364 controls in the training set. Among the recognized DMPs, 24,659 CpGs were hypermethylated and 15,356 CpGs were hypomethylated. The widespread distribution of DMPs on chromosomes hinted the crucial role of DNA methylation in the pathogenesis of CHD (Fig. 1A). In terms of quantity, the open sea region and the gene body were the main distribution regions of DMP (Fig. 1B, C). However, if the distribution density was considered, the CpG island and regions around transcriptional start site (TSS) were the most enriched regions. In addition, it was also observed that DMPs tended to appear in adjacent gene regions, such as TSS200 and 5'UTR, 5'UTR and 1stExon, gene body and 3'UTR (Fig. 1D). Considering that the regions surrounding TSS have received the most attention in CHD-related epigenetic researches, we further focused on the distribution of DMPs in TSS1500 and TSS200. Although DMPs were evenly spread in TSS1500 and TSS200 on each chromosome (Fig. 1E), certain differences were also observed in the distribution of DMPs in TSS1500 and TSS200. As shown in Fig. 1F, DMPs in TSS1500 were mainly located in CpG shores, while CpG islands were the principal distribution areas for DMPs in TSS200 (Fig. 1G).

We then performed gene set enrichment analysis (GSEA) based on the transcriptome data of the training set to probe into the underlying pathogenesis of CHD. The Gene Ontology (GO) analysis revealed that genes bounded up with autophagy of nucleus and cell redox homeostasis were up-regulated in CHD patients, while monocyte differentiation related genes were down-regulated (Fig. 1H). In Kyoto Encyclopedia of Genes and Genomes (KEGG) analysis, pathways involving fatty acid metabolism and lysosome were up-regulated in CHD (Fig. 1I). According to the results of Reactome analysis, genes correlated to autophagy and LDL clearance were notably up-regulated in CHD patients (Fig. 1J).

## WGCNA and core module analysis

Weighted correlation network analysis (WGCNA) was carried out in the training set to discriminate gene modules that linked with CHD related phenotypes (Fig. 2A). A total of 7 functional gene modules were recognized through WGCNA, among which the red module had the most prominent negative correlation with CHD (Fig. 2B). Besides, the red module was also negatively related to TG and positively related to HDL-C. Genes in the red module participated in a variety of functions and pathways that interrelated with CHD, including autophagosome assemble, cell cycle signaling, apoptosis, oxidative stress, and cholesterol metabolism (Fig. 2C–E).
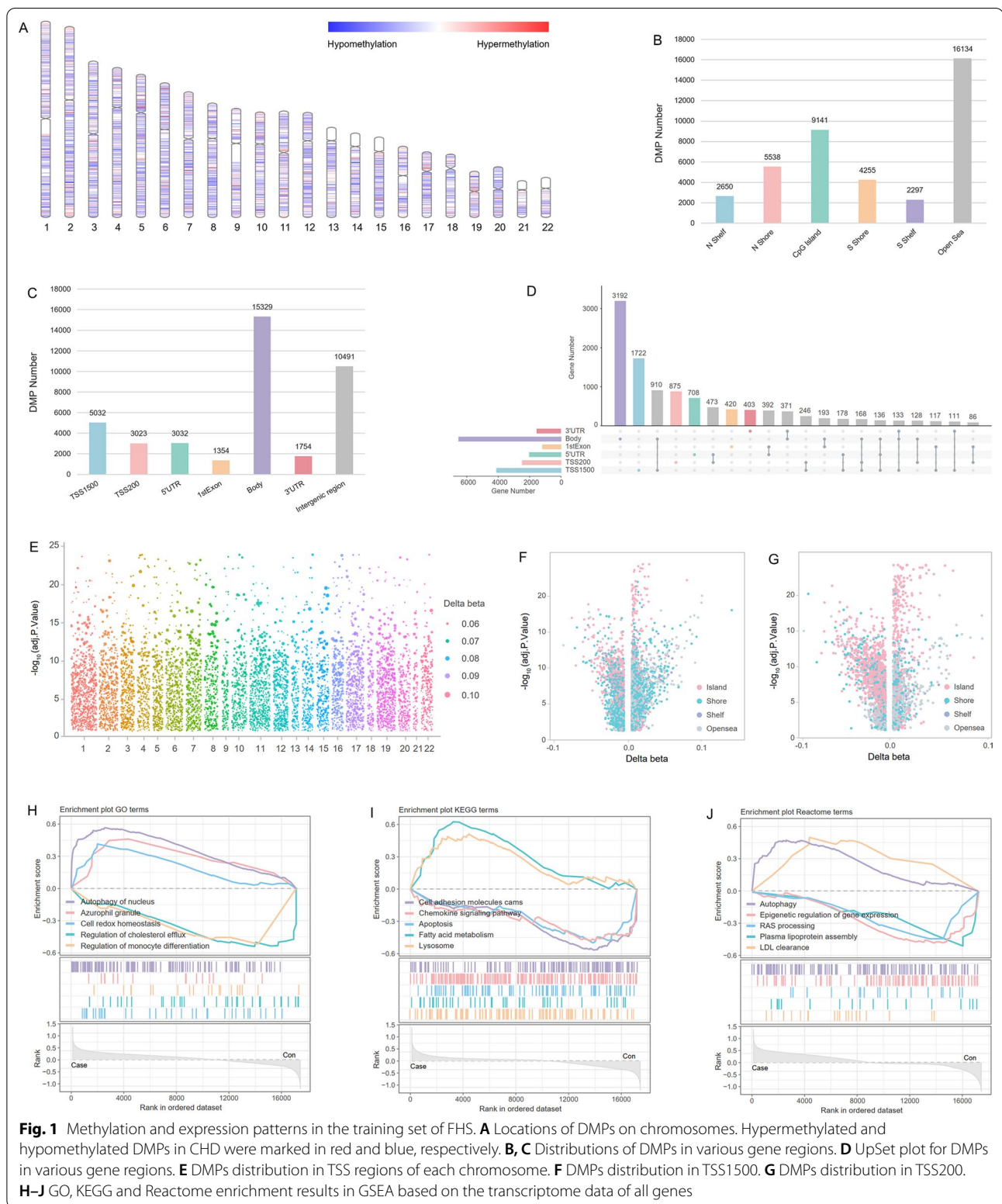
The red module was further visualized based on the connectivity obtained from the topological overlap matrix (TOM) (Fig. 2F). Core module of the TOM was consisted of the top 10 pivotal genes (Fig. 2G). Another network was established on the grounds of the protein interactions in the red module (Fig. 2H). The top 10 crucial genes of the protein–protein interaction (PPI) network are shown in Fig. 2I. It was noteworthy that the gene Autophagy Related 7 (ATG7) was the intersection of TOM core module and PPI network core module.

## Identification and function analysis of the hub genes

A total of 312 intersection genes between differentially methylated genes (DMGs) and differentially expressed genes (DEGs) were identified from the training set, among which 54 genes had significant correlations between methylation and expression statuses (Fig. 3A). These 54 genes were regarded as potential DNA methylation-regulated genes and the representative infor 4 mation were displayed in Fig. 3B. Similar to the results of WGCNA, the DNA methylation-regulated genes were also enriched in autophagy, oxidative stress, inflammation, immune cells function and proliferation related terms (Fig. 3C).

The methylation and expression statuses of the 54 DNA methylation-regulated genes were taken as 108 independent features and were all included into least absolute shrinkage and selection operator (LASSO) analysis. A total of 68 features were retained in the LASSO model (Fig. 3D, E). Of the retained features, only 13 genes had both methylation and expression features. In order to screen the features more rigorously, we introduced SelectFpr into the dimensionality reduction process. Finally, there were five genes selected by LASSO and SelectFpr simultaneously, namely BTB domain and CNC homolog 2 (BACH2), cyclin dependent kinase inhibitor 1B (CDKN1B), 24-dehydrocholesterol reductase (DHCR24), myeloperoxidase (MPO), and the aforementioned ATG7 (Fig. 3F).

The five DNA methylation-regulated genes mentioned above were incorporated into the subsequent machine learning model construction as hub genes of CHD. These hub genes were primarily enriched in autophagy, oxidative stress, lipoprotein metabolism, cell proliferation and aging (Fig. 3G). ATG7 ($p < 0.0001$), DHCR24 ($p < 0.0001$), MPO ($p < 0.0001$) were hypomethylated in CHD patients from the training set (Fig. 4A). Meanwhile, the expression levels of ATG7 ($p = 0.0238$), DHCR24 ($p = 0.0091$), MPO ($p = 0.0019$) were increased in CHD patients (Fig. 4B). Conversely, BACH2 ($p < 0.0001$) and CDKN1B ($p < 0.0001$) were hypermethylated in CHD patients. The

**Fig. 1** Methylation and expression patterns in the training set of FHS. **A** Locations of DMPs on chromosomes. Hypermethylated and hypomethylated DMPs in CHD were marked in red and blue, respectively. **B, C** Distributions of DMPs in various gene regions. **D** UpSet plot for DMPs in various gene regions. **E** DMPs distribution in TSS regions of each chromosome. **F** DMPs distribution in TSS1500. **G** DMPs distribution in TSS200. **H–J** GO, KEGG and Reactome enrichment results in GSEA based on the transcriptome data of all genes

Zhang *et al. Clinical Epigenetics*      (2022) 14:122

Page 5 of 17

expression levels of BACH2 ($p = 0.0421$) and CDKN1B ($p = 0.0001$) were declined in CHD patients. The expression levels of the five genes were all negatively correlated with the corresponding methylation status, whether in CHD patients, controls or all individuals (Fig. 4C–G).

### Machine learning modeling based on the hub genes

Decision curve analysis (DCA) and clinical impact curve analysis (CICA) were conducted to initially appraise the clinical practicability of the hub genes. We established three models in DCA and CICA based on the methylation parameters, expression parameters, and methylation-expression combination parameters of the five hub genes in the training set, respectively. The results from DCA indicated that when the risk threshold probability was set to CHD prevalence, which was 0.07 [17], the overall net benefit of the combination model was superior than unitary model based on methylation or expression parameters (Fig. 5A). Besides, CICA showed the false positive rate of the combination model was much lower than the methylation model and expression model when the risk threshold probability set to 0.07 (Fig. 5B).

Given that the appropriate algorithms could take full advantage of the current data to construct models with more excellent performance, we then established machine learning models using LightGBM, XGBoost and Random Forest. Finally, nine models were built through the three algorithms based on the methylation features, the expression features, and the methylation-expression combination features from the five hub genes. It could be observed that models formed by LightGBM always performed better than XGBoost and Random Forest models (Table 2). Compared with the methylation or expression models, the performances of the combination models were substantially improved (Fig. 5C–Q). Of all the nine models, the combination model established by Light-GBM was the top performer, whose F1 score and AUC were 0.517 and 0.834 in the validation set, respectively. In addition, the methylation and expression model established by LightGBM performed well in another two datasets obtained from the GEO database (Additional file 1: Table S1).

The traditional FRS model showed extremely poor capability in distinguishing CHD patients (Fig. 5R). The true positive rate (TPR), also known as sensitivity, was as low as 0.194 in the validation set of the FRS model. In contrast, the LightGBM-based combination model had considerable sensitivity in CHD identification, and its TPR reached 0.672 in the validation set. The true negative rate (TNR) of the LightGBM-based combination model even reached 0.864 in the validation set. In addition, the expression and methylation of ATG7 were identified as the dominant features in LightGBM-based models (Fig. 6A–F).

### Validation of the identified DNA methylation-regulated genes

Since monocytes were the principal cell type in peripheral blood involved in the pathogenesis of CHD, we verified the methylation and expression statuses of the five hub genes in monocytes. In accordance with the results from FHS, ATG7 ($p = 0.0080$), DHCR24 ($p = 0.0044$) and MPO ($p = 0.0465$) were hypomethylated in monocytes from CHD patients compared with controls (Fig. 6G). The expression levels of ATG7 ($p = 0.0146$) and DHCR24 ($p < 0.0001$) were up-regulated in CHD patients (Fig. 6H). Although no significant methylation difference was observed in BACH2, the expression of BACH2 was obviously reduced in CHD patients ($p = 0.0074$). The principal component analysis (PCA) revealed the five hub DNA methylation-regulated genes could effectively distinguish CHD patients from controls, especially when methylation and expression data were applied simultaneously (Fig. 6I–K).
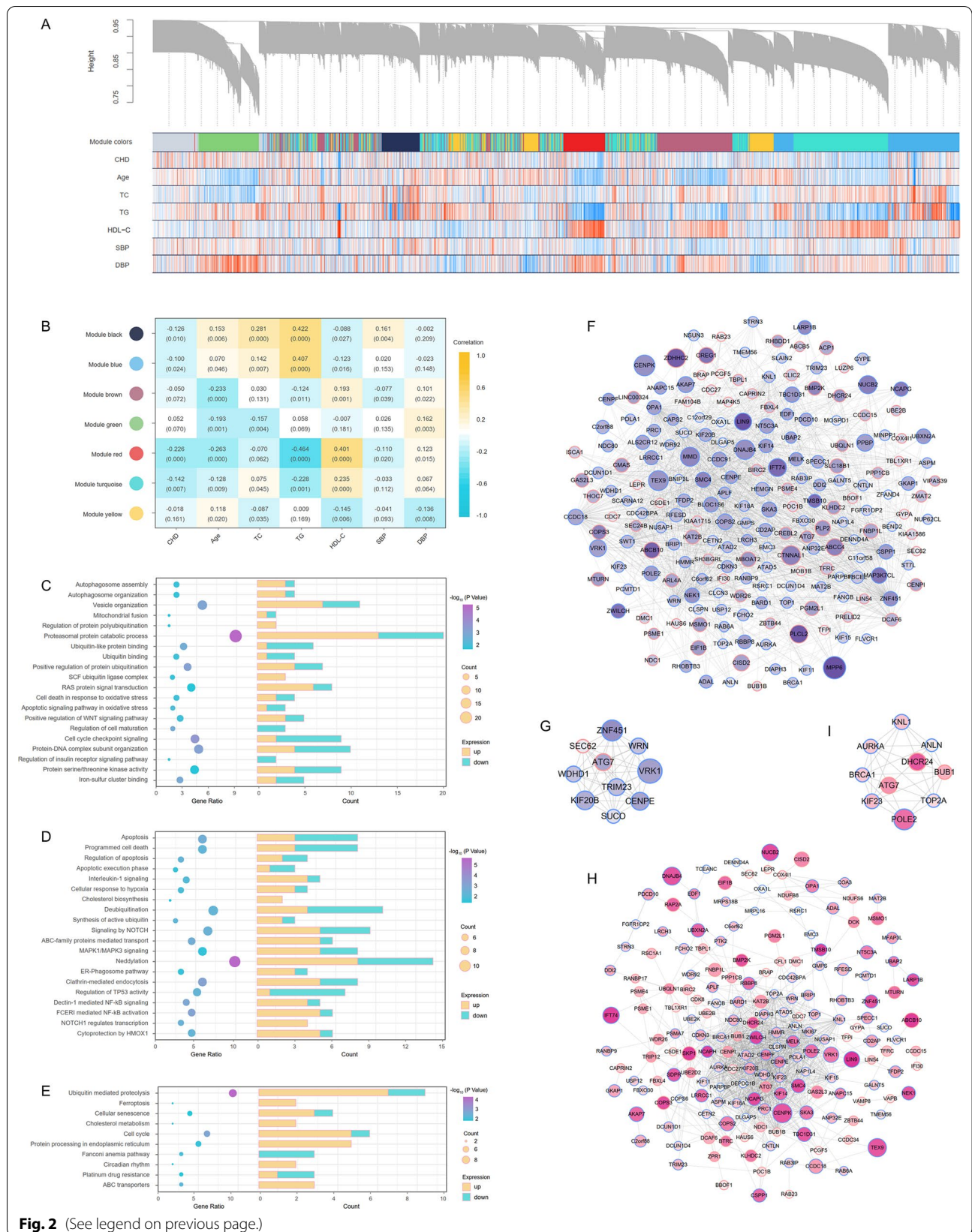
### Discussion

In the present study, the methylation and expression patterns in CHD were elucidated. Further, we constructed a machine learning model based on five identified DNA methylation-regulated genes, which exhibited eminent discriminability for CHD. The five DNA methylation-regulated genes were also validated in monocytes of our study population.

It has been extensively documented that DNA methylation is involved in the etiopathogenesis of CHD as a dominant epigenetic factor [11–13]. In addition to the previously reported enrichment of DMGs in TSS regions, CpG islands and adjacent intragenic regions [14, 15], we

---

(See figure on next page.)
**Fig. 2** WGCNA and module analysis based on the training set of FHS. **A** Modules identified from WGCNA were assigned unique colors. The heatmap showed the correlation between genes and phenotypes. Red and blue represented positive and negative correlations, respectively. **B** Correlations between modules and phenotypes. The top number was the correlation coefficient, and the bottom number was the corresponding *p* value. **C–E** GO (C), Reactome (D) and KEGG (E) enrichment results of the red module. **F, G** Network and core network based on connectivity of TOM in the red module. **H, I** Network and core network based on PPI in the red module. The diameter of each circle was positively correlated with the expression difference of the corresponding gene. Analogously, the deeper color indicated greater statistical significance. Red edge of circles meant the genes were up-regulated in CHD, while blue represented the genes were down-regulated in CHD

Zhang *et al. Clinical Epigenetics*     (2022) 14:122

Page 6 of 17



**Fig. 2** (See legend on previous page.)

noticed differences in the distribution of DMGs between TSS1500 and TSS200. There were totally 5032 and 3023 DMPs identified in TSS1500 and TSS200, respectively. In view of the linear length, it could be considered that the DMPs density of TSS200 was much higher than that of TSS1500. Besides, the DMPs located in the CpG shores accounted for 49.24% of all DMPs in TSS1500, but only occupied 18.82% of all DMPs in TSS200. On the contrary, DMPs in the CpG islands only made up 20.09% of all DMPs in TSS1500, but accounted for 51.70% of the DMPs in TSS200 located in CpG islands. The phenomenon not only showed the proportion of CpG island-shore was discrepant in TSS1500 and TSS200, but also hinted the difference of methylation pattern between TSS1500 and TSS200 in CHD.

A review written by Xia et al. summarized recent findings of large-scale epigenome-wide association studies (EWASs) related with incident CHD (iCHD) and CHD risk factors [5]. Similarly, our study identified five DNA methylation-regulated genes as epigenetic risk factors and elucidated the mechanisms of DNA methylation involved in the progression of CHD. GSEA showed that genes related to CHD were mainly enriched in autophagy, oxidative stress and monocyte regulation, highly consistent with the functions of the red module in WGCNA. Furthermore, ATG7 was identified as the hub gene of both the TOM-based network and the PPI-based network in the red module. ATG7 has been reported as an indispensable participant in the process of autophagy in atherosclerosis [16, 17]. The vital identity of ATG7 in the aforementioned two networks reconfirmed that autophagy played crucial roles in the pathological mechanism of CHD.

The enormous data from thousands of individuals in FHS were excellent resources for construction credible models via machine learning. The models based on the methylation features of DNA methylation-regulated genes performed much better than those using expression features. This superiority of models based on methylation data insinuated the DNA methylation statuses might have more potentials in CHD prediction than the mRNA expressions. The performance of models that constructed by methylation-expression combination features had made further leaps than the unitary methylation models. Although the specificity of the optimal combination model (0.864) was equivalent to that of the traditional FRS model (0.885), the sensitivity of the optimal combination model (0.672) was three times outnumbered the FRS model (0.194). Though it is desired for prediction tools to perform with high sensitivity and specificity, given the adverse effect of false negative, sensitivity should be considered as the more important metric in population screening [18].
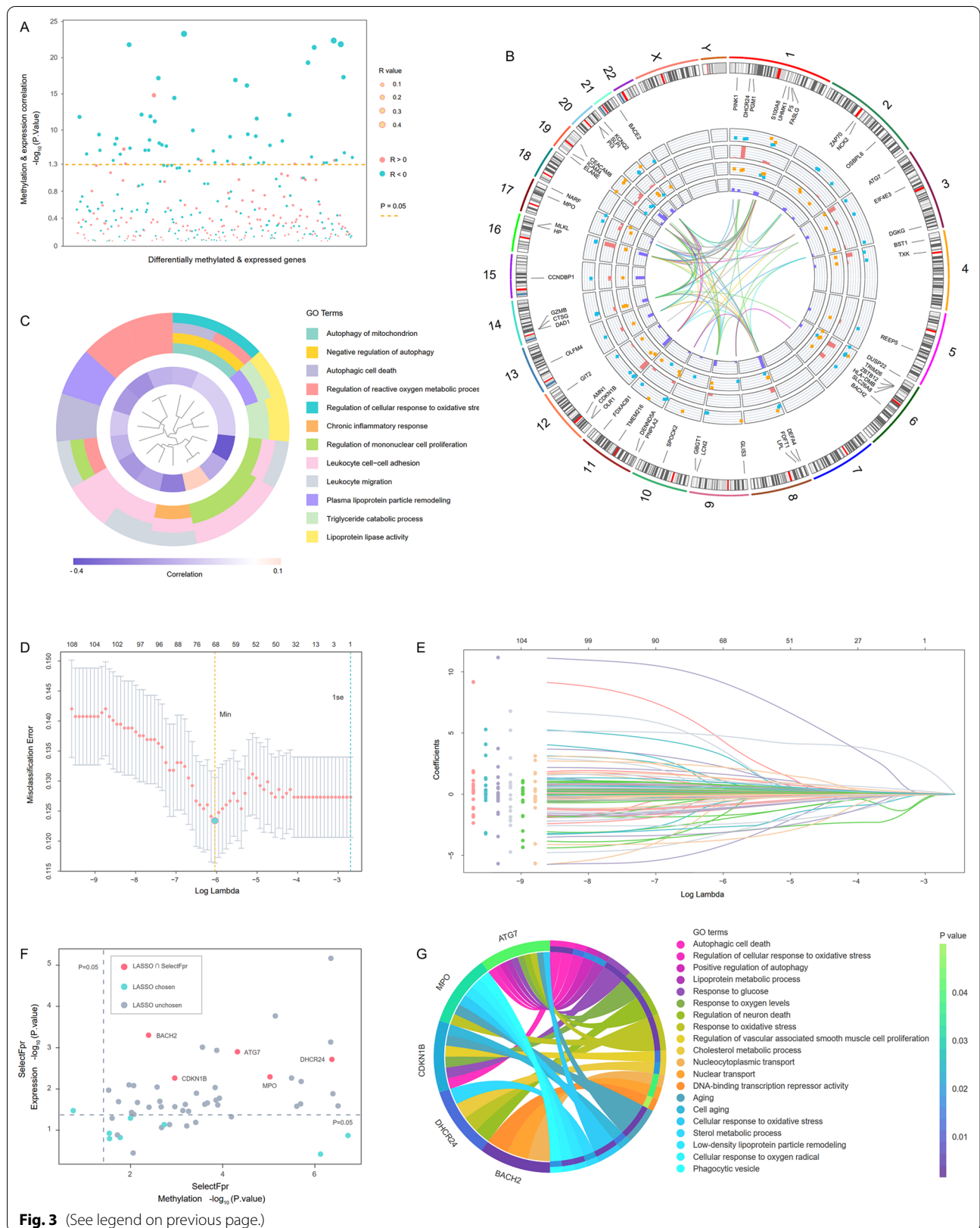
Navas-Acien et al. reported multiple DMPs associated with incident CHD in American Indians could be recurrent in other ethnicities, such as Hispanics, African Americans and Asians [19]. Similarly, of the five DNA methylation-regulated genes identified from PBLs in FHS, the differential methylation and/or expression of ATG7, BACH2, DHCR24 and MPO were further validated in monocytes from Chinese populations. The consistency confirmed the reliability of the bioinformatics analysis and suggested similarities of the methylation and expression patterns of these genes in the pathological process of CHD among multiple ethnicities. Even though no significant differences were observed for CDKN1B in monocytes, CDKN1B was still worthy of further study. The bias might be caused by the interference of other cell types in peripheral blood and the insufficient number of samples in the monocyte verification procedure.

Whether in the methylation model, expression model or combination model, the methylation and expression features of ATG7 always occupied decisive positions among all features, indicating ATG7 played an irreplaceable role in the predictive models of CHD. Integrating our results with published studies, the underlying pathways that ATG7 involved in the pathological progression of CHD were exhibited in Fig. 7. The hypomethylation in ATG7 promoter might enhance ATG7 expression in CHD. The up-regulated ATG7 catalyzed the formation of LC3-II by collaborating with LC3-I, ATG3 and phosphatidyl
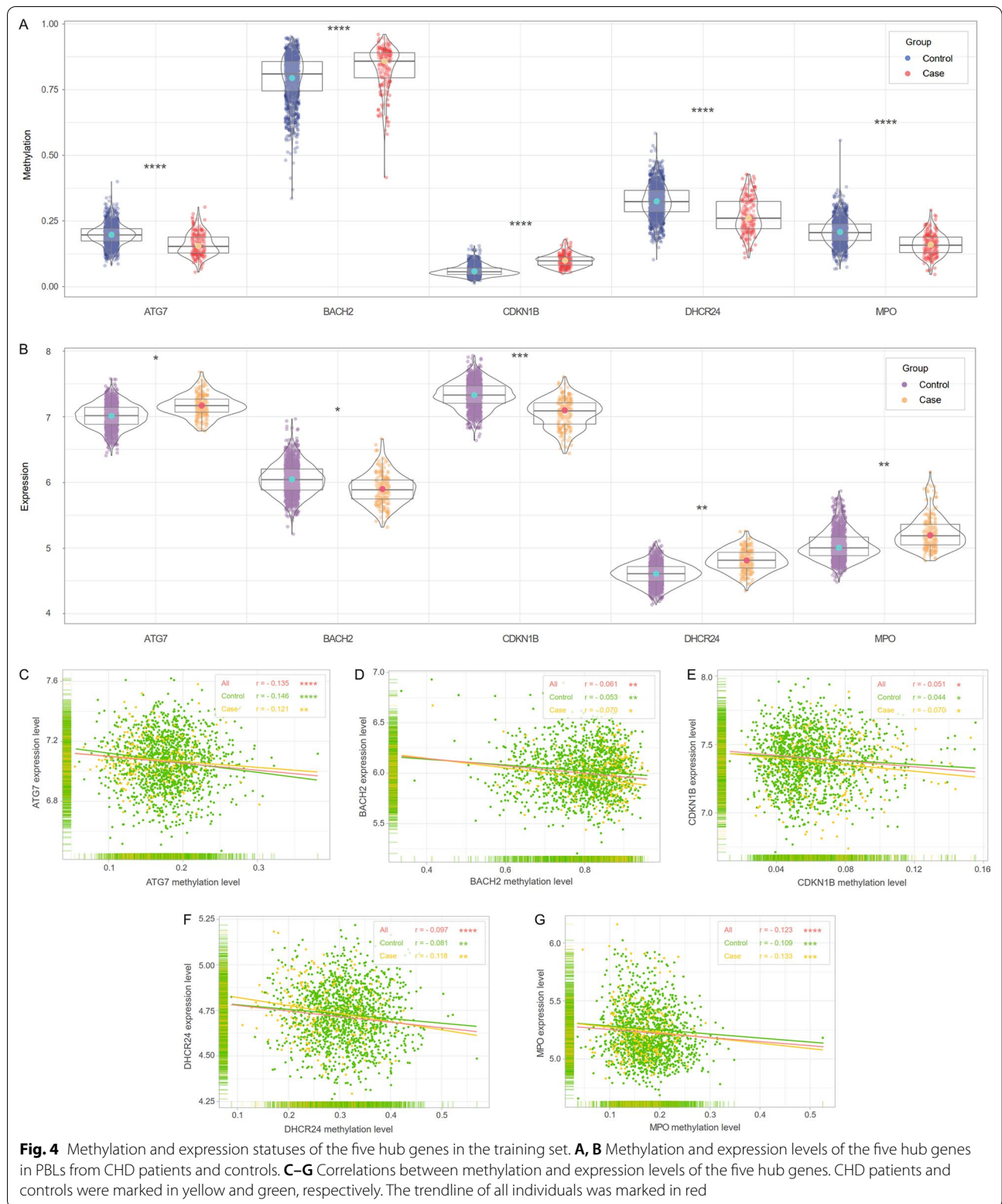
(See figure on next page.)
**Fig. 3** Dimensionality reduction to screen hub genes in the training set. **A** Part of the differentially methylated and expressed genes had statistically significant correlations between methylation and expression levels. The diameter of each circle was positively correlated with the absolute values of Spearman correlation coefficients. **B** Information of the potential DNA methylation-regulated genes. The circles from inside to outside were respectively: PPI networks among the genes, *p* value histogram and logFC scatter plot of the expression difference, *p* value histogram and delta beta scatter plot of the methylation difference, names and chromosomal locations of representative genes. **C** GO clustering enrichment results of the potential DNA methylation-regulated genes. The inner loop showed the Spearman correlation coefficients between methylation and expression status. **D** Tuning parameter selection of the potential DNA methylation-regulated genes by misclassification error in the LASSO model. **E** LASSO coefficients profiles of the features. **F** Intersection between LASSO and SelectFpr test. The dash lines indicated the threshold value of SelectFpr test. Genes chosen by both LASSO and SelectFpr were marked in red. **G** GO enrichment results of ATG7, BACH2, CDKN1B, DHCR24 and MPO

Zhang *et al. Clinical Epigenetics*      (2022) 14:122

Page 8 of 17



**Fig. 3** (See legend on previous page.)

Zhang *et al. Clinical Epigenetics*    (2022) 14:122

Page 9 of 17



**Fig. 4** Methylation and expression statuses of the five hub genes in the training set. **A, B** Methylation and expression levels of the five hub genes in PBLs from CHD patients and controls. **C–G** Correlations between methylation and expression levels of the five hub genes. CHD patients and controls were marked in yellow and green, respectively. The trendline of all individuals was marked in red

Zhang *et al. Clinical Epigenetics*     (2022) 14:122

Page 10 of 17

ethanolamine (PE) [20, 21]. The increased LC3-II facilitated the phagophore elongation [22], which might lead to monocyte dysfunction through aberrant excessive autophagy. DHCR24 was the pivotal gene that catalyzed the formation of cholesterol from desmosterol [23]. The up-regulation of DHCR24 might be induced by the hypomethylation of its promoter in CHD, leading to the immoderate consumption of desmosterol [24]. The diminished desmosterol in mitochondria allowed ROS production and thus promoting the glycolytic metabolic switch, a hallmark of inflammatory macrophages [24, 25]. The inhibitory effect of CKDN1B on macrophage proliferation might be reversed by hypermethylation of CKDN1B promoter [26]. Further, down-regulation of CDKN1B has been reported to induce atherosclerotic plaques formation and inflammatory response in the plaques. The transcription inhibition of BACH2 might be mediated by its promoter hypermethylation in CHD. The reduced BACH2 disabled the maintenance of its inhibitory effect on IFN-$\gamma$, resulting in IFN-$\gamma$ mediated inflammation in monocytes [27]. MPO released from monocytes was thought to contribute to endothelial dysfunction by limiting nitric oxide (NO) bioavailability via formation of reactive oxidants including hypochlorous acid (HOCl) [28]. The hypomethylation of MPO promoter in CHD might promoted MPO expression in monocyte, thereby accelerating vascular endothelial cell (VEC) injury. The dysfunction of monocytes caused by the above DNA methylation-regulated genes might affect the process of macrophages and foam cells formation, and then participated in the occurrence and development of CHD [29].

Some limitations in the present study deserved attentions. In view of FHS mainly contained European descents, the results in Chinese or other populations might deviate, though DNA methylation-regulated genes were verified in a small Chinese population. Moreover, the sample size was relatively insufficient for the validation of DNA methylation-regulated genes in monocytes. Besides, the epigenetic regulation mechanisms of ATG7 and other hub genes in CHD were also required to be further investigated in multi-ethnic populations.

## Conclusion

In summary, we identified five DNA methylation-regulated genes based on a predictive model for CHD using machine learning, which may clue the new epigenetic mechanism for CHD.

## Methods

### Methylome and transcriptome data in FHS

The dbGaP database authorized us to obtain FHS data (accession number: phs000007) and the Medical Ethics Committee of Zhongnan Hospital of Wuhan University approved all analyses described in this study (approval number: 2018017). The methylome and transcriptome data were collected from individuals of the Offspring cohort, who attended the eighth examination cycle of FHS. The CHD designation was evaluated and determined by a panel of three investigators from the Framingham Endpoint Review Committee. The basic criterion for CHD diagnosis was that at least one coronary artery had more than 50% stenosis. The CHD status of each individual was assessed in the eighth and ninth examination cycles from 2005 to 2008 and 2011 to 2014, respectively.
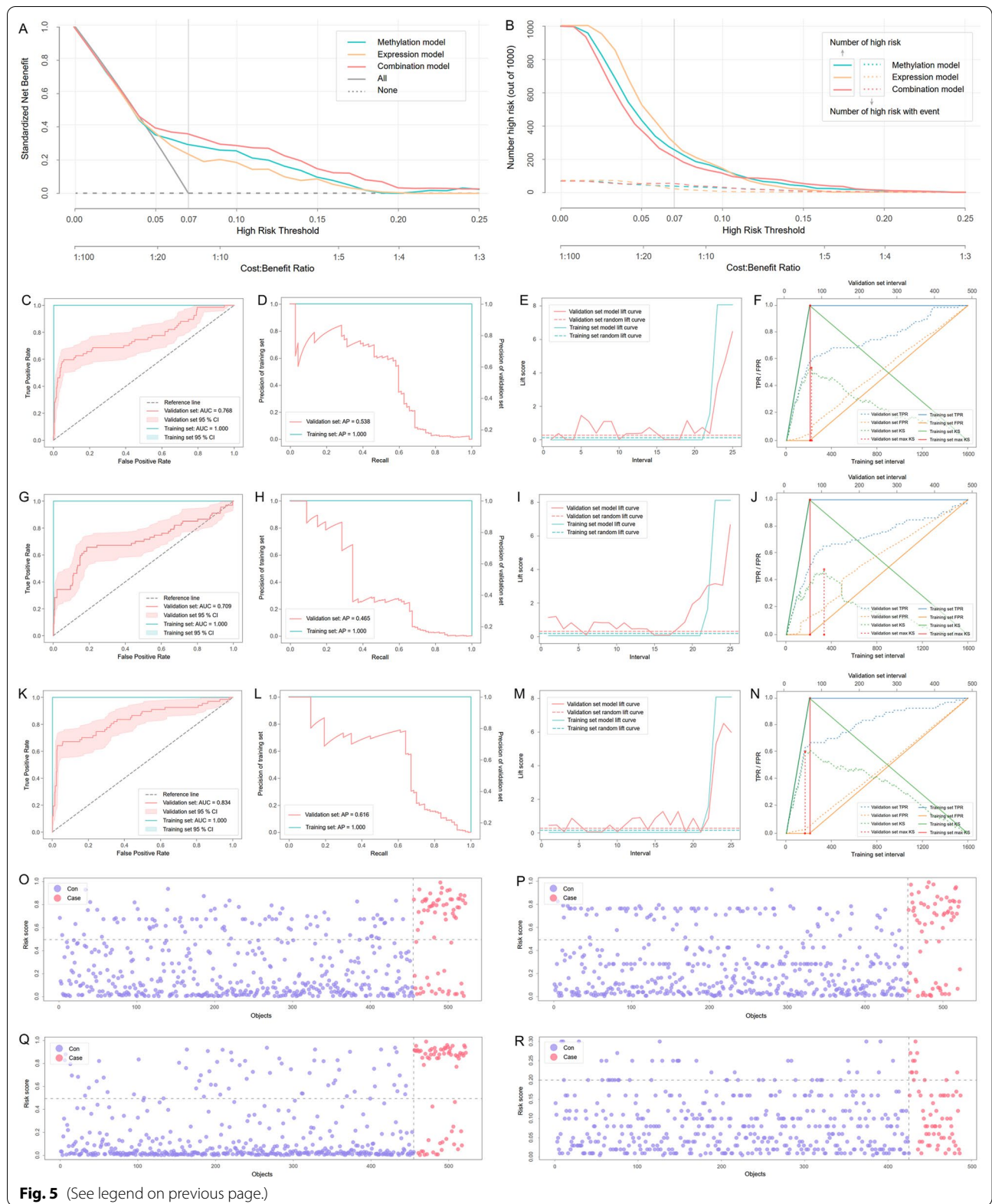
The DNA methylome data were obtained from PBLs of 2724 individuals and were detected by Illumina Infinium HumanMethylation450 BeadChip. The PBLs transcriptome data were available for 2441 individuals and were profiled by Affymetrix Human Exon 1.0 ST Array. There were a total of 2321 individuals provided both PBLs methylome and transcriptome data. In order to prevent the interference of age and gender between control group and CHD group, we adopted stratified sampling among the 2321 individuals to select age-gender-matched individuals. Through sampling, we selected 2117 age-gender-matched individuals ($p = 0.0700$ for age; $p = 0.0552$ for gender) to conduct subsequent quality control process.

### Quality control of methylome and transcriptome data

In the quality control process for methylome data, probes with detection $p > 0.01$ calculated from ChAMP package (version 2.21.1) were considered as failed probes and were excluded. Non-CpG probes, probes located on sexual chromosomes, and probes with bead count < 3 in at least 5% of the samples were also removed. Then, 6 samples were recognized as outliers by multi-dimensional

(See figure on next page.)

**Fig. 5** DCA, CICA and machine learning modeling. **A** DCA based on the five hub genes. The x-axis represented the threshold probability and the y-axis represented the net benefit. The gray line indicated a perfect prediction model, while the dash line represented a noneffective prediction model. The prevalence rate of CHD (0.07) was also marked in the figure. **B** CICA based on the five hub genes. The number of high-risk patients and the number of high-risk patients with events were drawn with solid lines and dash lines to represent different threshold probabilities, respectively. **C–N** Performance of the machine learning models constructed by LightGBM. ROC, PR, lift chart and KS plot of the methylation model (**C–F**), expression model (**G–J**) and combination model (**K–N**). **O–R** Risk plot of the validation set in the methylation model (**O**), expression model (**P**), combination model (**Q**) and FRS model (**R**)

**Fig. 5** (See legend on previous page.)

Zhang *et al. Clinical Epigenetics*      (2022) 14:122

Page 12 of 17

**Table 2** Performances of the models based on machine learning

| Features | Algorithm | Dataset | F1 | ACC | AUC (95% CI) | AP | KS | TP | FP | TN | FN | TPR | TNR | Kappa |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Methylation | LightGBM | Training | 0.995 | 0.999 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 197 | 0 | 1364 | 2 | 0.990 | 1.000 | 0.994 |
| | | Validation | 0.460 | 0.807 | 0.768 (0.694–0.843) | 0.538 | 0.540 | 43 | 77 | 378 | 24 | 0.642 | 0.831 | 0.353 |
| | XGBoost | Training | 1.000 | 1.000 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 199 | 0 | 1364 | 0 | 1.000 | 1.000 | 1.000 |
| | | Validation | 0.429 | 0.770 | 0.756 (0.683–0.830) | 0.391 | 0.525 | 45 | 98 | 357 | 22 | 0.672 | 0.785 | 0.308 |
| | Random forest | Training | 0.995 | 0.999 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 197 | 0 | 1364 | 2 | 0.990 | 1.000 | 0.994 |
| | | Validation | 0.443 | 0.803 | 0.737 (0.656–0.818) | 0.611 | 0.517 | 41 | 77 | 378 | 26 | 0.612 | 0.831 | 0.334 |
| Expression | LightGBM | Training | 0.992 | 0.998 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 196 | 0 | 1364 | 3 | 0.985 | 1.000 | 0.991 |
| | | Validation | 0.447 | 0.801 | 0.709 (0.626–0.792) | 0.465 | 0.472 | 42 | 79 | 376 | 25 | 0.627 | 0.826 | 0.337 |
| | XGBoost | Training | 0.997 | 0.999 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 198 | 0 | 1364 | 1 | 0.995 | 1.000 | 0.997 |
| | | Validation | 0.426 | 0.784 | 0.706 (0.646–0.766) | 0.538 | 0.494 | 42 | 88 | 367 | 25 | 0.627 | 0.807 | 0.309 |
| | Random forest | Training | 1.000 | 1.000 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 199 | 0 | 1364 | 0 | 1.000 | 1.000 | 1.000 |
| | | Validation | 0.283 | 0.592 | 0.647 (0.563–0.731) | 0.347 | 0.320 | 42 | 188 | 267 | 25 | 0.627 | 0.587 | 0.105 |
| Combination | LightGBM | Training | 1.000 | 1.000 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 199 | 0 | 1364 | 0 | 1.000 | 1.000 | 1.000 |
| | | Validation | 0.517 | 0.839 | 0.834 (0.770–0.897) | 0.616 | 0.615 | 45 | 62 | 393 | 22 | 0.672 | 0.864 | 0.427 |
| | XGBoost | Training | 1.000 | 1.000 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 199 | 0 | 1364 | 0 | 1.000 | 1.000 | 1.000 |
| | | Validation | 0.439 | 0.780 | 0.807 (0.740–0.874) | 0.460 | 0.566 | 45 | 93 | 362 | 22 | 0.672 | 0.796 | 0.322 |
| | Random forest | Training | 0.995 | 0.999 | 1.000 (1.000–1.000) | 1.000 | 1.000 | 197 | 0 | 1364 | 2 | 0.990 | 1.000 | 0.994 |
| | | Validation | 0.452 | 0.791 | 0.818 (0.758–0.878) | 0.599 | 0.487 | 45 | 87 | 368 | 22 | 0.672 | 0.809 | 0.340 |
| FRS | Framingham 10- | Training | – | 0.830 | 0.647 (0.606–0.687) | – | – | 39 | 105 | 1188 | 147 | 0.210 | 0.919 | – |
| | Year risk scale | Validation | – | 0.797 | 0.610 (0.536–0.684) | – | – | 12 | 49 | 376 | 50 | 0.194 | 0.885 | – |

*ACC* Accuracy, *AUC* Area under the receiver operating characteristic curve, *CI* Confidence interval, *AP* Average precision score, *KS* Kolmogorov–Smirnov, *TP* True positive, *FP* False positive, *TN* True negative, *FN* False negative, *TPR* True positive rate, *TNR* True negative rate, *FRS* Framingham risk score. Dashes meant the parameters were not applicable in the Framingham 10-year risk scale

scaling (MDS) plot and these samples were eliminated. The methylome data of the rest samples were normalized through PBC method to adjust the methylated and unmethylated intensities. We explored for batch effect using MDS plots, and the batch effect was controlled by the ComBat function of the ChAMP package. The RefBase function of the ChAMP package was used to obtain methylation-based estimates of the blood cell-type counts. And the cell heterogeneities among samples were further adjusted by the RefBase function.

As for the analysis of transcriptome data, probes without relevant gene names were removed, and the maximum values were taken in the circumstance that plural probes corresponded to the same gene name. We removed 26 samples as they were identified as outliers by MDS plot. The batch effect in the rest samples was adjusted by regressing out the batch variable with the ComBat function of the SVA package (version 3.42.0). We utilized the digital sorting algorithm package (DSA, version 1.0) to estimate the blood cell-type counts and adjust the cell heterogeneities of the expression data [30, 31].
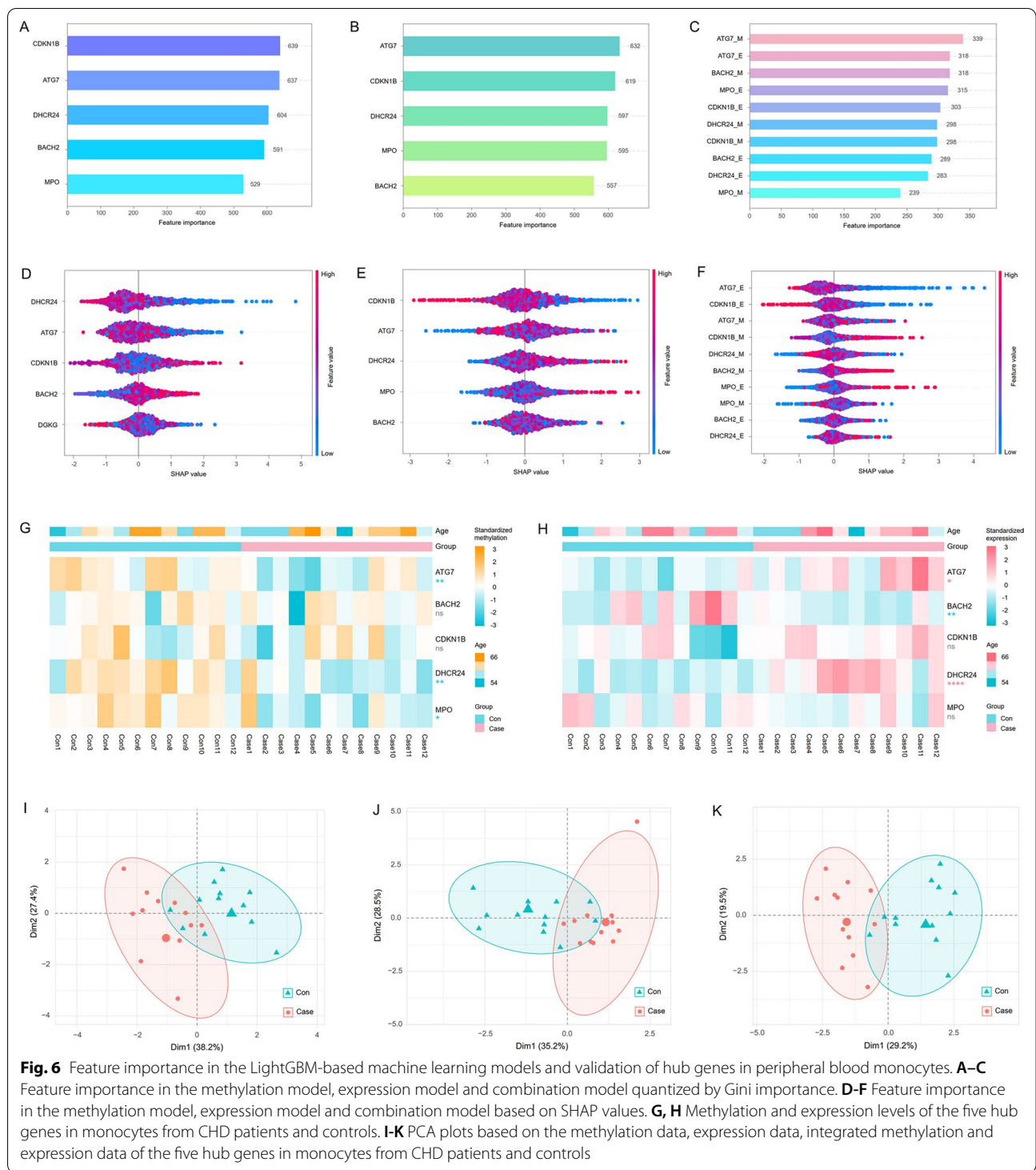
A total of 2085 samples passed the quality control process for methylation and expression data. These samples were randomly divided into training set and validation set. The training set contained 1563 individuals, of which 1364 individuals were controls and 199 individuals were CHD patients. The rest 522 individuals were classified into the validation set, of which 455 individuals were controls and 67 individuals were CHD patients. The clinical characteristics of the training set and the validation set were listed in Table 1. The subsequent bioinformatics analyses were based on the methylome and transcriptome data of the training set.

**Differential analysis of methylome and transcriptome data**

DMPs between CHD patients and controls in the training set were identified with threshold of *adj. p* < 0.05 and delta beta > 0.015. The *adj.p* was calculated by Benjamini and Hochberg (BH) method. Considering that the regions of 1500 and 200 bases upstream of the transcriptional start site (TSS1500; TSS200) are the paramount regulation regions [32–34], we took genes that had DMPs in the TSS1500 and TSS200 as DMGs. Further, the average beta values of DMPs in TSS1500 and TSS200 were regarded as the synthetic methylation statuses of the corresponding DMGs.

Empirical Bayes algorithm from Limma package (version 3.50.3) was utilized to recognize DEGs in the training set. Since using *adj. p* < 0.05 as threshold yielded very few DEGs, we chose *p* < 0.05 as threshold for DEGs screening.

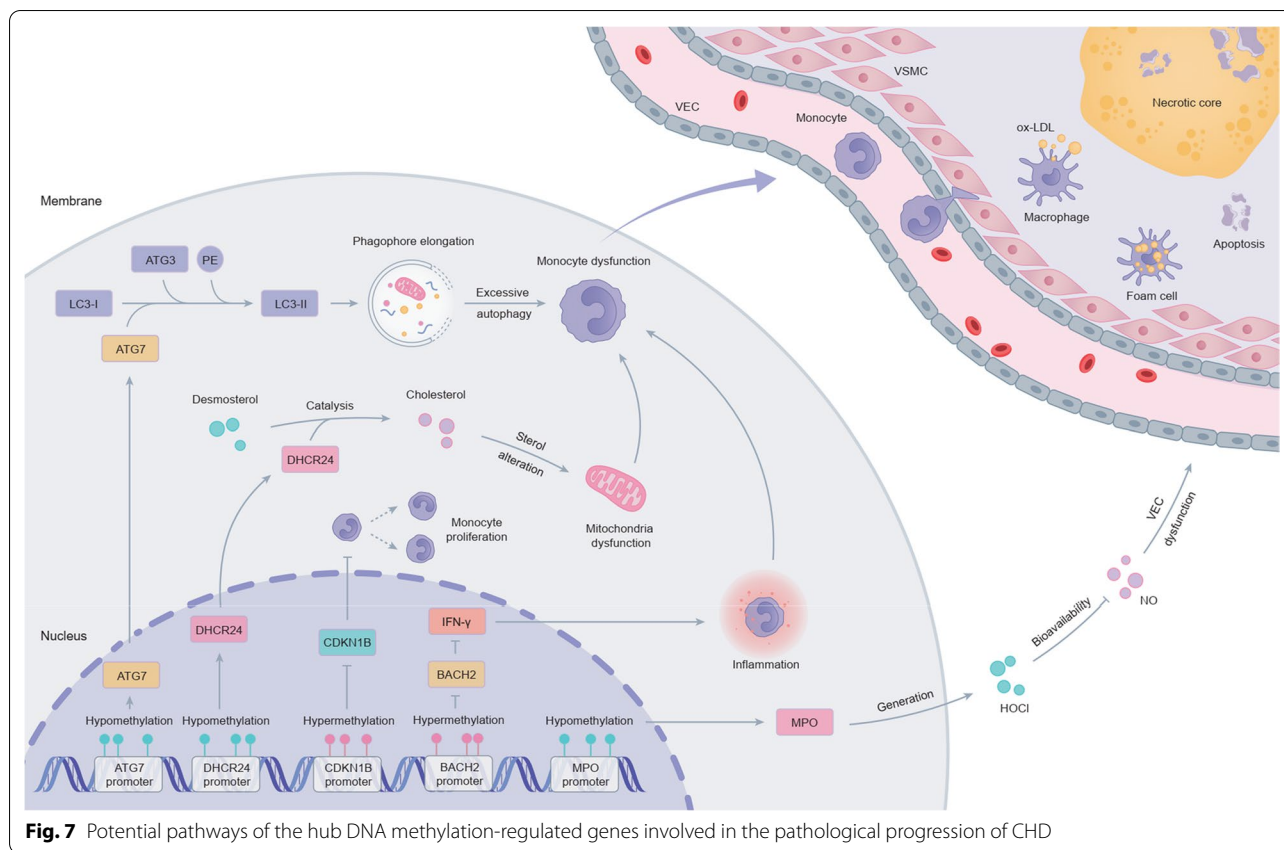Zhang *et al. Clinical Epigenetics*    (2022) 14:122

Page 13 of 17



**Fig. 6** Feature importance in the LightGBM-based machine learning models and validation of hub genes in peripheral blood monocytes. **A–C** Feature importance in the methylation model, expression model and combination model quantized by Gini importance. **D-F** Feature importance in the methylation model, expression model and combination model based on SHAP values. **G, H** Methylation and expression levels of the five hub genes in monocytes from CHD patients and controls. **I-K** PCA plots based on the methylation data, expression data, integrated methylation and expression data of the five hub genes in monocytes from CHD patients and controls

Further, among the intersection of DMGs and DEGs, genes with significant correlation between methylation and expression status were considered as DNA methylation-regulated genes. Spearman correlation test was applied in the evaluation of the correlation between methylation and expression status, and the threshold was $p < 0.05$.

**Enrichment analysis and WGCNA**

GO and KEGG enrichment analyses were performed using the ClusterProfiler package (version 4.2.2). While, Reactome enrichment analysis was carried out by ReactomePA package (version 1.38.0). The filter criteria for GO, KEGG and Reactome terms was $p < 0.05$. GSEA was

**Fig. 7** Potential pathways of the hub DNA methylation-regulated genes involved in the pathological progression of CHD

carried out by the GSEA software (version 4.2.1). We took c5.all.v7.2, c2.cp.kegg.v7.2 and c2.cp.reactome.v7.4 as the reference gene sets for GO, KEGG and Reactome in GSEA, respectively. Terms with normalized $p < 0.05$ were deemed to be statistically significant in GSEA.

Transcriptome data of the genes that in the top 25% of median absolute deviation (MAD) in the training set were used to conduct WGCNA by WGCNA package (version 1.71). During the sample clustering analysis, 43 outlying samples were eliminated when the clustering threshold was set to $-2.5$. The remaining 1520 samples were enrolled into the subsequent WGCNA procedure. With the soft threshold set to 5, the adjacency matrix was converted into TOM. Genes that were intrinsically homogeneous in expression patterns were divided into identical module. TOM-based connectivity network and PPI network were constructed in the module that was most intensely associated with CHD phenotype. The PPI network was established through Search Tool for the Retrieval of Interacting Genes (STRING) database with the interaction score set to 0.15. Core networks in the TOM network and PPI network were recognized based upon 12 algorithms provided by CytoHubba, a plug-in from Cytoscape (version 3.9.0).

**Dimensionality reduction process**

LASSO analysis has been widely applied in dimensionality reduction to select vital features for high-dimensional data [35]. The methylation and expression data of potential DNA methylation-regulated genes in the training set were imported into LASSO analysis as independent features by glmnet package (version 4.1.4). A gene was considered to be eligible only if both its methylation and expression features passed the screening of LASSO analysis.

Analogously, SelectFpr implements feature screening by performing false positive rate (FPR) estimation through scikit-learn (version 0.23.2). A gene was deemed to pass the SelectFpr test when its methylation and expression features were both with $p < 0.05$. The intersection genes between LASSO analysis and SelectFpr test were incorporated into subsequent machine learning modeling.

**DCA, CICA and machine learning modeling**

In order to preliminarily assess the clinical net benefits of the candidate DNA methylation-regulated genes in CHD prediction, we carried out DCA based on the methylation and expression data of the genes in the training set [36]. According to the epidemiological

Zhang *et al. Clinical Epigenetics*        (2022) 14:122

Page 15 of 17

statistics from 2006 to 2010 in USA, the prevalence of CHD was set to 0.07 in the DCA [37]. CICA was then performed on the basis of DCA to evaluate the practical value of the DNA methylation-regulated genes in CHD prediction. The DCA and CICA were both conducted by rmda package (version 1.6).

Machine learning modeling with appropriate algorithms can lead to observably potentiated predictive efficiency of biomarkers. In this study, LightGBM, XGBoost and Random Forest were utilized to construct machine learning models based on the hub DNA methylation-regulated genes, which were identified by the dimensionality reduction analyses. Each algorithm established 3 models using methylation data, expression data, and the combination of methylation and expression data, respectively.

The performance of models was assessed by a few parameters as following. The discernibility was evaluated by total accuracy (ACC), the classification accuracy was represented using Confusion matrix and Kappa value, while risk distinction and sorting ability was reflected by Kolmogorov–Smirnov (KS) and lift chart. Comprehensive performance of the models was estimated through receiver operating characteristic curve (ROC), area under ROC (AUC), precision-recall curve (PRC) and average precision score (AP). Finally, F1 score was considered as the paramount parameter that indicated the equilibrium of the models. The weight of features in the models was imputed by 2 methodologies, one was based on Gini importance, the other was in the light of SHapley Additive exPlanations (SHAP) values. The machine learning modeling was conducted on scikit-learn (version 0.23.2).

The optimal methylation model and expression model generated from FHS data were further applied to two additional datasets from Gene Expression Omnibus (GEO) database. GSE107143 contained peripheral blood methylome data from eight atherosclerosis patients and eight controls. GSE42148 provided peripheral blood transcriptome data from 13 CHD patients and 11 controls. The detection platforms of GSE107143 and GSE42148 were Illumina Infinium HumanMethylation450 BeadChip and Agilent SurePrint G3 Microarray, respectively. The data of these two datasets were processed and normalized with reference to the data processing procedures for FHS data.

## Wet experiment validation using transcriptome sequencing and methylation microarray

To further verify the identified DNA methylation-regulated genes, we performed experimental validation in Chinese populations. A total of 12 male CHD patients were recruited in Zhongnan Hospital of Wuhan University (Wuhan, China) between December 2019 and July 2021. The patients were diagnosed with CHD as coronary angiography confirmed that there were more that 50% stenotic lesions in at least one coronary artery. Another 12 males without history of cardiovascular events were recruited as controls. The clinical baseline information of the participants was listed in Additional file 1: Table S2.

We obtained 8 mL peripheral blood from each participant to isolate monocytes using Ficoll-Hypaque solution (Sigma-Aldrich, Germany) and CD14 microbeads (Miltenyi, Germany) according to the manufacturer instructions. Half of the monocytes were used for DNA extraction and the other half were used for RNA extraction (Omega, USA). An amount of 1 μg DNA was treated by bisulfate (Zymo Research, USA), followed by whole genome amplification, enzymatic end-point fragmentation, precipitation, and resuspension according to the Illumina Infinium HD Methylation Protocol. Then, the resuspended samples were hybridized on Illumina Infinium HumanMethylation850 BeadChip. Subsequent methylation analyses were conducted on ChAMP package (version 2.21.1). At the same time, over 1 μg total RNA was used to obtain mRNA by poly-T oligo-attached magnetic beads. The mRNA-seq libraries were established by the NEBNext Multiplex mRNA Library Prep Set for Illumina kit (New England Biolabs, USA) and finally sequenced on the Illumina NovaSeq 6000 platform. DESeq2 package (version 1.34.0) was applied to the sequencing data analyses.

## Statistical analysis

The Framingham 10-year risk scale is one of the most extensively used rating scales to evaluate the risk of developing cardiovascular diseases within ten years [38]. This scale assigns scores for a variety of traditional risk factors, including gender, age, TC, smoking, HDL-C and SBP [39]. The accumulated score of the 6 factors is referred to as FRS, which corresponds to the probability of suffering cardiovascular diseases in ten years. Individuals with estimated probabilities $\geq 20\%$ are considered as high-risk populations [40]. Among the 1563 individuals in the training set, parameters required for FRS calculation were applicable in 1479 individuals. Besides, 487 of the 522 individuals in the validation set provided parameters required for FRS calculation. We calculated the FRS using the Framingham 10-year risk scale and compared the performance between the FRS model and the DNA methylation-regulated genes models.

In the present study, normally distributed data were shown as mean $\pm$ standard deviation (SD), while data with non-normal distribution were described

Zhang *et al. Clinical Epigenetics*      (2022) 14:122

Page 16 of 17

as median and inter-quartile range. The differences between continuous data were analyzed by Student's $t$ test or Mann–Whitney $U$ test. Analyses involving binary data were processed by Chi-square test or Fisher's exact test. Pearson and Spearman correlation tests were applied to assess the correlations between continuous data. All statistical analyses were performed on GraphPad Prism (version 9.0), R (version 4.1.3) and Python (version 2.7.11). The threshold for the statistical tests was two-tail $p < 0.05$.

## Abbreviations
BMI: Body mass index; CHD: Coronary heart disease; CICA: Clinical impact curve analysis; DBP: Diastolic blood pressure; DCA: Decision curve analysis; DEG: Differentially expressed gene; DMG: Differentially methylated gene; DMP: Differentially methylated position; FBG: Fasting blood glucose; FHS: Framingham Heart Study; FRS: Framingham risk score; GO: Gene ontology; GSEA: Gene set enrichment analysis; HDL-C: High-density lipoprotein cholesterol; KEGG: Kyoto Encyclopedia of Genes and Genomes; LASSO: Least absolute shrinkage and selection operator; LDL-C: Low-density lipoprotein cholesterol; PBL: Peripheral blood leukocyte; PPI: Protein–protein interaction; SBP: Systolic blood pressure; SHAP: SHapley additive explanations; TC: Total cholesterol; TG: Triglycerides; TOM: Topological overlap matrix; TSS: Transcriptional start site; WGCNA: Weighted correlation network analysis.

## Supplementary Information
The online version contains supplementary material available at https://doi.org/10.1186/s13148-022-01343-2.

> **Additional file 1: Table S1**. Performance of the expression and methylation models applied to datasets from GEO. **Table S2**. Clinical baseline information of participants in the transcriptome sequencing and methylation microarray.

## Declarations

### Ethics approval and consent to participate
This study was approved by Medical Ethics Committee of Zhongnan Hospital of Wuhan University approved all analyses described in this study (approval number: 2018017).

### Consent for publication
Not applicable.

### Competing interests
The authors declare that they have no competing interests.

### Author details
[1]Center for Gene Diagnosis and Department of Clinical Laboratory Medicine, Zhongnan Hospital of Wuhan University, Donghu Road 169, Wuhan 430071, China. [2]Department of Clinical Laboratory Medicine, Renmin Hospital of Wuhan University, Wuhan 430060, China.

## References
1. Tsao CW, et al. Heart disease and stroke statistics-2022 update: a report from the American Heart Association. Circulation. 2022;145:e153–639. https://doi.org/10.1161/cir.0000000000001052.
2. Dogan MV, Grumbach IM, Michaelson JJ, Philibert RA. Integrated genetic and epigenetic prediction of coronary heart disease in the Framingham Heart Study. PLoS ONE. 2018;13: e0190549. https://doi.org/10.1371/journal.pone.0190549.
3. Wang C, et al. Development and validation of a predictive model for coronary artery disease using machine learning. Front Cardiovasc Med. 2021;8: 614204. https://doi.org/10.3389/fcvm.2021.614204.
4. Dogan MV, Knight S, Dogan TK, Knowlton KU, Philibert R. External validation of integrated genetic-epigenetic biomarkers for predicting incident coronary heart disease. Epigenomics. 2021;13:1095–112. https://doi.org/10.2217/epi-2021-0123.
5. Xia Y, Brewer A, Bell JT. DNA methylation signatures of incident coronary heart disease: findings from epigenome-wide association studies. Clin Epigenetics. 2021;13:186. https://doi.org/10.1186/s13148-021-01175-6.
6. Xie M, et al. BMAL1-downregulation aggravates porphyromonas gingivalis-induced atherosclerosis by encouraging oxidative stress. Circ Res. 2020;126:e15–29. https://doi.org/10.1161/circresaha.119.315502.
7. Branco MR, et al. Maternal DNA methylation regulates early trophoblast development. Dev Cell. 2016;36:152–63. https://doi.org/10.1016/j.devcel.2015.12.027.
8. Khera R, et al. Use of machine learning models to predict death after acute myocardial infarction. JAMA Cardiol. 2021;6:633–41. https://doi.org/10.1001/jamacardio.2021.0122.
9. Peng X, Li L, Wang X, Zhang H. A machine learning-based prediction model for acute kidney injury in patients with congestive heart failure. Front Cardiovasc Med. 2022;9: 842873. https://doi.org/10.3389/fcvm.2022.842873.
10. Chen G, Li L, Tao H. Bioinformatics identification of ferroptosis-related biomarkers and therapeutic compounds in ischemic stroke. Front Neurol. 2021;12: 745240. https://doi.org/10.3389/fneur.2021.745240.
11. Dai Y, Chen D, Xu T. DNA methylation aberrant in atherosclerosis. Front Pharmacol. 2022;13: 815977. https://doi.org/10.3389/fphar.2022.815977.
12. Bakshi C, Vijayvergiya R, Dhawan V. Aberrant DNA methylation of M1-macrophage genes in coronary artery disease. Sci Rep. 2019;9:1429. https://doi.org/10.1038/s41598-018-38040-1.
13. Do WL, et al. Associations between DNA methylation and BMI vary by metabolic health status: a potential link to disparate cardiovascular outcomes. Clin Epigenet. 2021;13:230. https://doi.org/10.1186/s13148-021-01194-3.
14. Zhang X, et al. Identification of potential biomarkers for CAD using integrated expression and methylation data. Front Genet. 2020;11:778. https://doi.org/10.3389/fgene.2020.00778.
15. Zhang X, et al. DYSF promotes monocyte activation in atherosclerotic cardiovascular disease as a DNA methylation-driven gene. Trans Res J Lab Clin Med. 2022. https://doi.org/10.1016/j.trsl.2022.04.001.
16. Osonoi Y, et al. Defective autophagy in vascular smooth muscle cells enhances cell death and atherosclerosis. Autophagy. 2018;14:1991–2006. https://doi.org/10.1080/15548627.2018.1501132.

Zhang *et al. Clinical Epigenetics*     (2022) 14:122

Page 17 of 17

17. Amersfoort J, et al. Defective autophagy in T cells impairs the development of diet-induced hepatic steatosis and atherosclerosis. Front Immunol. 2018;9:2937. https://doi.org/10.3389/fimmu.2018.02937.

18. Dogan MV, et al. Blood-based biomarkers for predicting the risk for five-year incident coronary heart disease in the Framingham Heart Study via machine learning. Genes. 2018. https://doi.org/10.3390/genes9120641.

19. Navas-Acien A, et al. Blood DNA methylation and incident coronary heart disease: evidence from the Strong Heart Study. JAMA Cardiol. 2021;6:1237–46. https://doi.org/10.1001/jamacardio.2021.2704.

20. Naponelli V, Modernelli A, Bettuzzi S, Rizzi F. Roles of autophagy induced by natural compounds in prostate cancer. Biomed Res Int. 2015;2015:121826. https://doi.org/10.1155/2015/121826.

21. Huang F, Mai J, Chen J, He Y, Chen X. Non-coding RNAs modulate autophagy in myocardial ischemia-reperfusion injury: a systematic review. J Cardiothorac Surg. 2021;16:140. https://doi.org/10.1186/s13019-021-01524-9.

22. Sureshbabu A, Ryter SW, Choi ME. Oxidative stress and autophagy: crucial modulators of kidney injury. Redox Biol. 2015;4:208–14. https://doi.org/10.1016/j.redox.2015.01.001.

23. Capell-Hattam IM, et al. Twin enzymes, divergent control: the cholesterogenic enzymes DHCR14 and LBR are differentially regulated transcriptionally and post-translationally. J Biol Chem. 2020;295:2850–65. https://doi.org/10.1074/jbc.RA119.011323.

24. Zhang X, et al. Desmosterol suppresses macrophage inflammasome activation and protects against vascular inflammation and atherosclerosis. Proc Nat Acad Sci USA. 2021. https://doi.org/10.1073/pnas.2107682118.

25. Nishimura T, et al. Hepatitis C virus impairs p53 via persistent overexpression of 3beta-hydroxysterol delta24-reductase. J Biol Chem. 2009;284:36442–52. https://doi.org/10.1074/jbc.M109.043232.

26. Yamada S, et al. Inhibition of local macrophage growth ameliorates focal inflammation and suppresses atherosclerosis. Arterioscler Thromb Vasc Biol. 2018;38:994–1006. https://doi.org/10.1161/atvbaha.117.310320.

27. Okkenhaug K, Graupera M, Vanhaesebroeck B. Targeting PI3K in cancer: impact on tumor cells, their protective stroma, angiogenesis, and immunotherapy. Cancer Discov. 2016;6:1090–105. https://doi.org/10.1158/2159-8290.Cd-16-0716.

28. Cheng D, et al. Inhibition of MPO (myeloperoxidase) attenuates endothelial dysfunction in mouse models of vascular inflammation and atherosclerosis. Arterioscler Thromb Vasc Biol. 2019;39:1448–57. https://doi.org/10.1161/atvbaha.119.312725.

29. Luo Y, et al. Macrophagic CD146 promotes foam cell formation and retention during atherosclerosis. Cell Res. 2017;27:352–72. https://doi.org/10.1038/cr.2017.8.

30. Zhong Y, Wan YW, Pang K, Chow LM, Liu Z. Digital sorting of complex tissues for cell type-specific gene expression profiles. BMC Bioinform. 2013;14:89. https://doi.org/10.1186/1471-2105-14-89.

31. Tcw J, et al. Cholesterol and matrisome pathways dysregulated in astrocytes and microglia. Cell. 2022;185:2213-2233.e2225. https://doi.org/10.1016/j.cell.2022.05.017.

32. Luo R, et al. DNA methylation subpatterns at distinct regulatory regions in human early embryos. Open Biol. 2018. https://doi.org/10.1098/rsob.180131.

33. Lu Z, et al. Locus-specific DNA methylation of Mecp2 promoter leads to autism-like phenotypes in mice. Cell Death Dis. 2020;11:85. https://doi.org/10.1038/s41419-020-2290-x.

34. Yagi M, et al. Identification of distinct loci for de novo DNA methylation by DNMT3A and DNMT3B during mammalian development. Nat Commun. 2020;11:3199. https://doi.org/10.1038/s41467-020-16989-w.

35. Mullah MAS, Hanley JA, Benedetti A. LASSO type penalized spline regression for binary data. BMC Med Res Methodol. 2021;21:83. https://doi.org/10.1186/s12874-021-01234-9.

36. Kerr KF, Brown MD, Zhu K, Janes H. Assessing the clinical impact of risk prediction models with decision curves: guidance for correct interpretation and appropriate use. J Clinic Oncol. 2016;34:2534–40. https://doi.org/10.1200/jco.2015.65.5654.

37. Prevalence of coronary heart disease--United States, 2006–2010. MMWR. Morbidity and mortality weekly report. (2011);**60**,1377–1381 .

38. Xu L, et al. Viscus fat area contributes to the Framingham 10-year general cardiovascular disease risk in patients with type 2 diabetes mellitus. Life Sci. 2019;220:69–75. https://doi.org/10.1016/j.lfs.2019.01.036.

39. Executive summary of the third report of the National cholesterol education program (NCEP) expert panel on detection, evaluation, and treatment of high blood cholesterol in adults (adult treatment panel III). Jama. 2001; **285**, 2486–2497. https://doi.org/10.1001/jama.285.19.2486.

40. D'Agostino RB Sr, Grundy S, Sullivan LM, Wilson P. Validation of the Framingham coronary heart disease prediction scores: results of a multiple ethnic groups investigation. JAMA. 2001;286:180–7. https://doi.org/10.1001/jama.286.2.180.

## Publisher's Note