

1 **Title:** Machine Learning Reveals the Contribution of Rare Genetic Variants and  
2 Enhances Risk Prediction for Coronary Artery Disease in the Japanese Population

3

4 **Authors:**

5 Hirotaka Ieki <sup>1,2,3,4</sup>, Kaoru Ito <sup>1\*</sup>, Sai Zhang <sup>5</sup>, Satoshi Koyama <sup>1,6,7</sup>, Martin Kjellberg <sup>3,8</sup>,

6 Hiroki Yoshida <sup>1,2</sup>, Ryo Kurosawa <sup>1</sup>, Hiroshi Matsunaga <sup>1,2</sup>, Kazuo Miyazawa <sup>1</sup>,

7 Nobuyuki Enzan <sup>1,6,7</sup>, Changhoon Kim <sup>9</sup>, Jeong-Sun Seo <sup>9,10</sup>, Koichiro Higasa <sup>11,12</sup>,

8 Kouichi Ozaki <sup>1,13</sup>, Yoshihiro Onouchi <sup>1,14</sup>, The Biobank Japan Project, Koichi Matsuda <sup>15</sup>,

9 Yoichiro Kamatani <sup>16</sup>, Chikashi Terao <sup>17</sup>, Fumihiko Matsuda <sup>12</sup>, Michael Snyder <sup>3,4\*</sup>,

10 Issei Komuro <sup>18,19\*</sup>

11

12 **Affiliations:**

13 <sup>1</sup> Laboratory for Cardiovascular Genomics and Informatics, RIKEN Center for  
14 Integrative Medical Sciences, Yokohama, Japan.

15 <sup>2</sup> Department of Cardiovascular Medicine, Graduate School of Medicine, The  
16 University of Tokyo, Tokyo, Japan.

17 <sup>3</sup> Department of Genetics, Center for Genomics and Personalized Medicine, Stanford  
18 University School of Medicine, Stanford, USA.

19 <sup>4</sup> Stanford Cardiovascular Institute, Stanford University School of Medicine, Stanford,  
20 USA.

21 <sup>5</sup> Department of Epidemiology, University of Florida, Gainesville, USA.

22 <sup>6</sup> Center for Genomic Medicine, Department of Medicine, Massachusetts General  
23 Hospital, Boston, MA, USA.

24 <sup>7</sup> Program in Medical and Population Genetics, Broad Institute of Harvard and MIT,

25 Cambridge, MA, USA.

26 <sup>8</sup> School of Engineering Sciences in Chemistry, Biotechnology and Health, Royal

27 Institute of Technology (KTH), Stockholm, Sweden

28 <sup>9</sup> Bioinformatics Institute, Macrogen Inc., Seoul, Republic of Korea.

29 <sup>10</sup> Asian Genome Institute, Seoul National University Bundang Hospital, Gyeonggi-do,

30 Republic of Korea.

31 <sup>11</sup> Department of Genome Analysis, Institute of Biomedical Science, Kansai Medical

32 University, Hirakata, Japan.

33 <sup>12</sup> Human Disease Genomics, Center for Genomic Medicine, Kyoto University

34 Graduate School of Medicine, Kyoto, Japan.

35 <sup>13</sup> Medical Genome Center, Research Institute, National Center for Geriatrics and

36 Gerontology, Obu, Japan.

37 <sup>14</sup> Department of Public Health, Chiba University Graduate School of Medicine, Chiba,

38 Japan.

39 <sup>15</sup> Department of Computational Biology and Medical Science, Graduate School of

40 Frontier Sciences, The University of Tokyo, Tokyo, Japan.

41 <sup>16</sup> Laboratory of Complex Trait Genomics, Department of Computational Biology and

42 Medical Sciences, Graduate School of Frontier Sciences, The University of Tokyo,

43 Tokyo, Japan.

44 <sup>17</sup> Laboratory for Statistical and Translational Genetics, RIKEN Center for Integrative

45 Medical Science, Yokohama, Japan

46 <sup>18</sup> International University of Health and Welfare, Tokyo, Japan

47 <sup>19</sup> Department of Frontier Cardiovascular Science, Graduate School of Medicine, The

48 University of Tokyo, Tokyo, Japan

49

50

51 **Correspondence:**

52 kaoru.ito@riken.jp (K.I.)

53 mpsnyder@stanford.edu (M.S.)

54 komuro-ky@umin.ac.jp (I.K.)

55

56 **Summary**

57 Genome-wide association studies (GWASs) have advanced our understanding of  
58 coronary artery disease (CAD) genetics and enabled the development of polygenic risk  
59 scores (PRSs) for estimating genetic risk based on common variant burden. However,  
60 GWASs have limitations in analyzing rare variants due to insufficient statistical power,  
61 thereby constraining PRS performance. Here, we conducted whole genome sequencing  
62 of 1,752 Japanese CAD patients and 3,019 controls, applying a machine learning-based  
63 rare variant analytic framework. This approach identified 59 CAD-related genes,  
64 including known causal genes like *LDLR* and those not previously captured by GWASs.  
65 A rare variant-based risk score (RVS) derived from the framework significantly  
66 predicted CAD cases and cardiovascular mortality in an independent cohort. Notably,  
67 combining the RVS with traditional PRS improved CAD prediction compared to PRS  
68 alone (area under the curve, 0.66 vs 0.61;  $p=0.007$ ). Our analyses reinforce the value of  
69 incorporating rare variant information, highlighting the potential for more  
70 comprehensive genetic assessment.

71

72 **Keywords**

73 Coronary artery disease; Rare variants; Genetic risk estimation; Polygenic risk  
74 score; Myocardial infarction; Genome-wide association study; Machine learning;  
75 BioBank Japan

76

77

78

## 79 **Introduction**

80           Despite advancements in treatments and medications, coronary artery disease  
81 (CAD), encompassing conditions such as angina pectoris and myocardial infarction  
82 (MI), remains a leading cause of death worldwide <sup>1,2</sup>. CAD etiology is complex,  
83 involving a multifaceted interplay between genetic predisposition and environmental  
84 determinants. Lifestyle factors including diet, smoking, and physical activity are well-  
85 established contributors to the onset and progression of CAD <sup>3,4</sup>. Additionally,  
86 conditions such as elevated low-density lipoprotein (LDL) cholesterol, hypertension,  
87 and glucose intolerance further exacerbate the risk profile <sup>5</sup>. The importance of genetic  
88 predisposition is also underscored by a European twin study, which estimated that  
89 genetic factors contributed to over 50% of CAD development <sup>6,7</sup>. Therefore,  
90 understanding the genetic underpinnings of CAD and accurately estimating an  
91 individual's lifetime genetic risk are crucial for effective prevention and management  
92 strategies.

93           To date, genome-wide association studies (GWASs) and their meta-analyses  
94 have identified more than 300 loci associated with CAD <sup>8-12</sup>. Polygenic risk scores  
95 (PRSs) derived from GWAS summary statistics have enabled the estimation of  
96 individual-level CAD risk <sup>13,14</sup>. However, despite these significant advancements, the  
97 heritability of CAD explained by GWASs remains lower than anticipated. This gap  
98 may be partly attributed to the primary focus of GWAS on low frequency to common  
99 variants, while rare variants are often underrepresented in these analyses <sup>5,15</sup>. Rare  
100 variants often have a large effect size on diseases and phenotypes, making them a  
101 promising target for drug development <sup>16</sup>. Incorporating rare variants into genetic risk  
102 scores could significantly enhance the accuracy of CAD prediction. Despite this

103 potential, previous GWASs and aggregated rare variant association analyses have  
104 struggled even in large-scale sequencing studies, identifying only a few genes at exome-  
105 wide significance per trait <sup>17,18</sup>. Furthermore, calculating a genetic risk score based on  
106 rare variants is challenging because gene-level effect sizes are not estimated by  
107 conventional gene-based analysis methods.

108           Recently, advancements in machine learning have led to the development of  
109 novel methods for genetic analysis, one of which is the HEAL (Hierarchical Estimate  
110 from Agnostic Learning) method, a machine learning-based framework for  
111 comprehensive rare variant analysis. This approach has been successful in identifying  
112 disease-associated genes and creating genetic risk scores in patients with abdominal  
113 aortic aneurysm <sup>19</sup>. In the current study, we conducted whole genome sequencing  
114 (WGS) of Japanese CAD patients and applied a modified version of the HEAL  
115 framework tailored for CAD to analyze rare variants and systematically prioritize  
116 disease-associated genes. Furthermore, we developed a rare variant-based genetic risk  
117 score (RVS) using this framework and validated the performance with an independent  
118 cohort. We then explored the relationship between the RVS and GWAS-based PRS to  
119 elucidate the characteristics of rare variants in CAD, bridging the gap in our  
120 understanding of CAD genetics by incorporating rare variant information, potentially  
121 uncovering novel insights into disease mechanisms and improving risk prediction  
122 models.

## 123 **Results**

### 124 **Whole genome sequencing of CAD samples in the Japanese population**

125 The overview and the design of our study are shown in **Figure 1**. We  
126 performed WGS on the discovery cohort comprising 1,765 Japanese CAD patients and  
127 3,148 controls. In order to enhance the genetic discovery power<sup>20</sup>, we prioritized  
128 patients with early-onset MI, a severe form of CAD, from the BioBank Japan (BBJ)  
129 cohort. The average age of MI onset in these patients was  $47.4 \pm 4.1$  years, indicating a  
130 relatively young population with a severe disease phenotype. After quality control of the  
131 WGS data, we retained 4,771 individuals (1752 cases and 3019 controls) with  
132 51,717,580 genetic variants. For the validation WGS cohort, we included 200 CAD  
133 cases and 824 control samples with 25,531,471 variants (**Table S1 and S2**).  
134 Demographic features in each cohort are summarized in **Table 1**. We then used the  
135 quality-controlled data for further analyses including single variant association tests to  
136 identify individual variants associated with CAD, a conventional gene-based association  
137 test to examine the cumulative effect of variants within specific genes, and a machine  
138 learning-based framework to uncover the potential contribution of rare variants (**Figure**  
139 **S1**).

140 We first conducted a single variant association test in the discovery cohort  
141 using a logistic regression model implemented in PLINK software with covariates of  
142 age, sex and top ten ancestry principal components (PCs). The genomic inflation factor  
143  $\lambda_{GC}$  was calculated to be 1.03, indicating minimal inflation of test statistics and  
144 suggesting that the quality control applied to the samples was adequate (**Figure S2**).  
145 This initial single variant association analysis did not identify any genetic loci that  
146 reached a genome-wide significance threshold of  $P = 5 * 10^{-8}$ . A subsequent analysis

147 was performed using SAIGE software designed to handle both common and rare  
148 variants, adjusted for age, sex and top ten ancestry PCs. This analysis revealed two  
149 previously reported loci on chromosome 12 that reached a genome-wide significance  
150 threshold (rs7977233;  $p=1.47 * 10^{-8}$ , rs3782886;  $p=1.47 * 10^{-8}$ , respectively, (**Figure S3**  
151 **and Table S3**))<sup>10,11</sup>. However, these were both common variants, emphasizing the  
152 difficulty in analyzing rare variants using current GWAS approaches.

153 To increase the detection power of rare variant associations, gene-based tests  
154 are often used, in which variants are aggregated and analyzed together for each gene.  
155 This approach allows for the analysis of rare variants that are underpowered in single  
156 variant association tests due to their low frequency. It also increases detection power by  
157 reducing the multiple testing burden. Thus, we conducted a gene-based rare variant  
158 aggregated association analysis using the sequential kernel association test-optimal  
159 (SKAT-O). While no genomic inflation was observed ( $\lambda = 0.939$ ) (**Figure S4**), the  
160 *LDLR* gene surpassed a suggestive threshold ( $p = 2.3 \times 10^{-5}$ ). However, no genes reached  
161 the gene-wide significance threshold of  $p = 2.5 \times 10^{-6}$  (**Figure S4 and Table S4**). This  
162 result also highlighted the challenges of analyzing rare variants in genetic association  
163 studies due to insufficient statistical power with a limited sample size.

164

## 165 **The machine learning-based framework prioritizes disease-associated genes and** 166 **reveals molecular networks**

167 We next conducted a machine learning-based rare variant analysis using a  
168 modified HEAL<sup>19</sup>. In this framework, we first quantified the mutation burden for each  
169 gene in each participant defined by the cumulative effects of deleterious  
170 nonsynonymous variants within the gene. We then trained a penalized logistic



171 regression model to predict disease status based on these mutation burden scores. The  
172 model was trained to identify a minimal set of most distinguishing features (genes) for  
173 CAD, while also optimizing parameters for accurate disease prediction. Through robust  
174 cross-validation (**Figure S5**), we successfully prioritized fifty-nine candidate genes  
175 associated with CAD development (**Table S5, S6 and Figure S6**).

176 To investigate the functions of the fifty-nine HEAL<sub>CAD</sub> genes, we assessed  
177 constraint scores and checked for overlaps with neighboring genes identified in previous  
178 GWASs on CAD and its risk factors. Using the Genehancer database <sup>21</sup>, which provides  
179 information on genome-wide enhancers and their target genes, we identified prioritized  
180 genes that overlapped with the target genes of enhancers found significant in previous  
181 GWASs. We also referenced the International Mouse Phenotyping Consortium (IMPC)  
182 <sup>22</sup> database to investigate the phenotypes associated with a gene knockout (KO) in mice  
183 and conducted gene set enrichment analysis to identify functional clusters among the  
184 HEAL<sub>CAD</sub> genes. The genes were subsequently categorized into eight distinct clusters  
185 based on the hierarchical clustering of their functional annotations (**Figure 2A, 2B and**  
186 **Table S7**).

187 Among these clusters, cluster 3 notably included the *LDLR* gene, which  
188 exhibited the strongest contribution to CAD. *LDLR* is a well-established causal gene for  
189 familial hypercholesterolemia <sup>23</sup> and has been consistently associated with CAD in  
190 previous GWASs and genome sequencing studies <sup>9,24,25</sup>, supporting the validity of our  
191 machine learning-based framework. In the IMPC database, *LDLR* KO mice showed  
192 increased circulating cholesterol levels <sup>26</sup>, a known risk factor for CAD. Cluster 7  
193 contained genes related to obesity and metabolic processes, such as the *RNF216* locus,  
194 which is associated with body mass index (BMI) <sup>27</sup> and increased glucose levels in KO

195 mice<sup>22</sup>. Additionally, the *VRK2* locus has been reported to be associated with BMI<sup>28</sup>,  
196 smoking behavior and alcohol use<sup>29</sup>, indicating its broader impact on metabolic health.  
197 Cluster 2 comprised genes identified by previous GWAS on phenotypes such as blood  
198 pressure, diabetes, and cholesterol levels. The *FTO* gene within this cluster was  
199 highlighted for its strong association with obesity<sup>30,31</sup> and related phenotypes linked to  
200 BMI<sup>32</sup>, LDL cholesterol<sup>33</sup>, blood pressure<sup>34</sup>, and CAD<sup>35</sup>. Cluster 8 encompassed  
201 genes associated with cholesterol levels, obesity and blood pressure in GWAS and  
202 GeneHancer categories, with phenotypic evidence in human and KO mice. For instance,  
203 the *CYP27A1* locus is associated with diastolic blood pressure<sup>36</sup> and triglyceride levels  
204<sup>37</sup> and has connections to cholesterol levels and premature CAD according to human  
205 phenotype ontology<sup>38</sup>.

206 To further determine the functions of the fifty-nine genes, we mapped them  
207 onto the human protein-protein interaction (PPI) network followed by identifying  
208 proteins that were tightly clustered with these HEAL<sub>CAD</sub> genes as topological modules  
209<sup>19</sup>. We identified 46 tightly clustered topological modules encompassing the HEAL<sub>CAD</sub>  
210 genes. Gene ontology analysis confirmed the functional coherence of the proteins  
211 within each module, revealing significant enrichment for specific biological processes.  
212 For instance, module M119 was significantly enriched for lipid homeostasis with a false  
213 discovery rate (FDR) of  $2.53 \times 10^{-22}$ , suggesting a critical role in regulating lipid levels  
214 (**Figure 2C and Table S8**). These modules included pathways known as CAD risk  
215 factors, such as lipid and glucose metabolism (M25, M31, M51, M86, M119). Notably,  
216 M119 included lipid metabolism-related genes such as *LDLR*, *PCSK9*, *LIPA*, and  
217 *ANGPTL3* (**Figure 2D**), which are well-known targets for medications treating  
218 dyslipidemia and CAD<sup>39 40</sup>. Other modules were associated with different biological

219 processes, including platelet volume (e.g., M13), immune system function (M1), blood  
220 vessel and heart development (e.g., M47, M328), and RNA metabolism and translation  
221 processes (e.g., M3, M34). While recent studies have indicated the contribution of  
222 common variants identified by CAD-GWAS to the disease through various pathways  
223 such as plaque formation, inflammation, transcriptional regulation, and angiogenesis<sup>41</sup>,  
224 our findings suggest that diverse biological processes are also implicated in CAD, even  
225 in the context of rare variants. This underscores the complexity of CAD pathogenesis,  
226 involving a wide array of biological pathways and molecular mechanisms.

227

#### 228 **Rare variant risk-based risk score and its clinical impact**

229 In conjunction with the prioritization of disease-related genes, the modified  
230 HEAL enabled us to develop a prediction model for CAD based on genetic information.  
231 Using the optimized machine learning model, we computed a rare variant-based risk  
232 score (RVS) for each individual. The RVS demonstrated a significant predictive  
233 capability for CAD, with an area under the receiver operating characteristics curve  
234 (AUROC) of 0.574, as validated through a nested cross-validation approach in the  
235 discovery cohort. When applied to an independent validation cohort, the RVS also  
236 identified CAD cases with an AUROC of 0.581 ( $p = 0.002$ ), indicating its ability to  
237 discriminate CAD cases.

238 To further understand the characteristics of RVS in terms of clinical aspects,  
239 we explored the association of RVS with clinically relevant parameters. The RVS  
240 showed significant correlations with several key clinical measurements, including low-  
241 density-lipoprotein cholesterol (LDLC), total bilirubin (TBil), alanine aminotransferase  
242 (ALT), prothrombin time (PT-INR), total cholesterol levels, neutrophil count, and

243 potassium levels (**Figure 3A and Table S9**). These correlations are noteworthy since  
244 elevated cholesterol levels and coagulation abnormalities are established risk factors for  
245 CAD<sup>42-44</sup>. Moreover, alterations of total bilirubin and AST were also reported to be  
246 associated with cardiovascular risk<sup>45,46</sup>, reinforcing the clinical relevance of the RVS in  
247 the context of CAD.

248 We extended our analysis to assess the impact of the RVS on long-term  
249 cardiovascular mortality. In the validation cohort, a higher RVS was significantly  
250 associated with increased cardiovascular mortality ( $P = 0.01$ , log-rank test) (**Figure 3B**).  
251 When exclusively analyzing CAD patients, those with higher RVS also exhibited a  
252 significantly worse cardiovascular mortality rate ( $p = 0.03$ , log-rank test) (**Figure 3C**).  
253 These findings suggest that RVS not only predicts CAD occurrence but also correlates  
254 with the disease severity and its long-term prognosis, highlighting its potential clinical  
255 utility in risk stratification and prognosis estimation for CAD patients.

256

### 257 **The integration of RVS and PRS improves the performance of the genomic risk** 258 **score**

259 Many GWASs have been conducted for CAD, leading to the development of  
260 PRS that primarily comprise common variants to predict the risk of CAD. Multiple  
261 studies have reported that PRS can serve as an important indicator for predicting and  
262 assessing the severity of CAD. Whereas these scores typically focus on common  
263 variants and do not account for rare variants, which can also significantly contribute to  
264 disease risk, our RVS encompasses rare variants not included in PRS. Thus, to compare  
265 the properties between RVS and PRS, we first calculated individual PRS based on  
266 CAD-GWAS<sup>11</sup> in the validation cohort. The PRS also significantly predicted CAD with

267 an AUROC of 0.61 ( $p = 0.001$ ; 95% confidence interval (C.I.), 0.565-0.653).  
268 Interestingly, there was no significant correlation between PRS and RVS ( $r = -0.01$ ,  $p =$   
269 0.73) (**Figure 4A**), indicating that RVS provides a different genomic perspective on  
270 CAD risk.

271 When examining CAD cases specifically, RVS showed a negative correlation  
272 with PRS ( $r = -0.17$ ,  $p = 0.015$ ) (**Figure 4A**). Additionally, PRS was associated with  
273 different clinical measurements compared to RVS, such as triglycerides, uric acid, body  
274 mass index (BMI), and activated partial thromboplastin time (APTT) and it was  
275 negatively associated with HDL cholesterol (HDLC), which is considered protective  
276 against CAD (**Figure 3A, Figure S7 and Table S10**)<sup>47</sup>. These data support the notion  
277 that PRS and RVS may have complementary rather than redundant roles in predicting  
278 CAD, as they were associated with different clinical parameters and did not show a  
279 positive correlation.

280 Given these distinct properties, we integrated PRS and RVS to develop a  
281 combined risk score (CRS) aiming at enhancement of the performance of the  
282 framework in predicting CAD. The CRS showed positive correlations with several  
283 clinical measures, including serum uric acid, coagulation functions, LDLC, and  
284 triglycerides (TG), while negatively correlating with HDLC levels (**Figure 4B and**  
285 **Table S11**). Focusing on lipid metrics, CRS demonstrated correlations with LDLC, TC,  
286 TG, and HDLC, suggesting that it combines the unique predictive elements of both RVS  
287 and PRS (**Figure 4C**). Finally, we evaluated the predictive performance of CRS and  
288 observed a significant improvement in CAD prediction compared to PRS alone in the  
289 validation cohort (AUROC 0.66 vs 0.61,  $p=0.007$ ; Pseudo  $R^2$  0.093 vs 0.040,  $p =$   
290 0.0018; AUPRC 0.35 vs 0.29,  $p = 0.0154$ ) (**Figure 5 and Table S12**). These results

291 suggest that RVS can complement PRS and that incorporating rare variant information  
292 as an RVS into PRS significantly enhances the ability to predict CAD, thereby  
293 addressing some of the unexplained heritability in the disease.  
294  
295

## 296 Discussion

297 In this study, we developed a machine learning-based analytical framework to  
298 investigate the genetics of CAD pathogenesis with a focus on rare variants. We  
299 leveraged this framework together with whole-genome sequencing (WGS) data from  
300 the Japanese population to enhance our understanding of the complex CAD genetic  
301 architecture. Our findings indicated that the modified HEAL, a machine learning-based  
302 framework, effectively prioritized genes associated with CAD, including the well-  
303 established *LDLR* gene, while also uncovering intricate molecular networks involved in  
304 the disease. The rare variant-based risk score (RVS) generated through this framework  
305 demonstrated significant predictive power for CAD and long-term cardiovascular  
306 mortality. Furthermore, the RVS showed different characteristics from conventional  
307 common variant-based PRS, and combining the rare variant-based RVS with the PRS  
308 substantially improved CAD prediction.

309 Identifying disease-associated rare variants remains a significant challenge,  
310 not only in single variant association analyses but also in aggregated rare variant  
311 association analyses<sup>48,49</sup>. While some studies have adopted a targeted resequencing  
312 approach by selecting specific genes based on prior knowledge<sup>25,50</sup>; previous attempts  
313 at genome-wide or exome-wide analyses have often suffered from insufficient statistical  
314 power, leading to limited success in identifying previously uncharacterized genes  
315 associated with complex traits like CAD<sup>20</sup>. Also in this study, the single variant  
316 association analysis and the gene-based rare variant association analysis failed to reveal  
317 genome-wide significant rare variants linked to CAD. Even in previous studies  
318 involving more than 450,000 exome sequencing data from the UK biobank, only a  
319 single gene, *LDLR*, reached a significance level in the gene-based test for CAD<sup>17</sup>.

320 These persistent challenges highlight the difficulties in rare variant analyses.

321 To address these challenges, we utilized a machine learning-based framework  
322 to analyze rare variants, building on the HEAL model in a prior study, where Li et al.  
323 successfully uncovered the genetic architecture of rare variants in abdominal aortic  
324 aneurysm<sup>19</sup>. We adapted and optimized the model for CAD patients, marking the first  
325 application of the technique in this disease context. Unlike the previous HEAL model  
326 that focused only on missense single nucleotide variants (SNVs), our approach casts a  
327 wider net as it incorporates insertion, deletion and putative loss-of-function (pLOF)  
328 variants. This comprehensive inclusion of variant types allows for a more holistic  
329 examination of the genetic landscape underlying CAD, potentially capturing a broader  
330 spectrum of disease-associated genetic alterations. Furthermore, the robustness of our  
331 model was enhanced by hyperparameter tuning through a grid search to avoid  
332 overfitting and we evaluated its predictive performance using both internal cross-  
333 validation and an independent validation cohort<sup>51</sup>.

334 Through this improved framework, we successfully prioritized CAD-  
335 associated genes, extending beyond previously reported genes such as *LDLR*, *FTO*, and  
336 *CYP27A1*. By mapping these genes onto the human protein-protein interaction network,  
337 we uncovered 46 tightly clustered topological modules, providing insights into their  
338 functional roles in CAD pathogenesis. Beyond lipid metabolism, the analysis revealed  
339 modules associated with other relevant biological processes, including platelet function,  
340 immune system regulation, blood vessel and heart development, and RNA metabolism.  
341 Interestingly, while previous GWASs have highlighted the role of common variants in  
342 CAD development through various pathways, our findings suggest that rare variants  
343 also contribute to the disease through a wide spectrum of biological processes.



344           We also utilized our framework to develop an RVS and demonstrated its  
345   discriminative capacity between CAD cases and controls in the validation cohort. The  
346   distinctive feature of RVS lies in its utilization of rare nonsynonymous variants as input  
347   data, setting it apart from conventional PRS that primarily focus on common variants.  
348   This approach allows RVS to tap into a different spectrum of genomic information,  
349   involving risk factors uncaptured by PRS. The independence of RVS from PRS is  
350   further substantiated by the absence of a significant positive correlation between these  
351   two scoring systems and the complementary relationships with clinical risk parameters.  
352   This lack of correlation suggests that the RVS and PRS are capturing distinct aspects of  
353   genetic risk for CAD, each contributing unique information to the overall risk  
354   assessment. Importantly, the integration of RVS and PRS resulted in improved  
355   predictive performance, demonstrating a synergistic effect that enhanced the ability to  
356   accurately assess CAD risk. While methods combining information from one or a few  
357   genetic mutations with PRS have been reported <sup>52</sup>, our study presented a more  
358   comprehensive approach to combine rare and common variant information.  
359   Furthermore, these findings reinforce the recognition that rare variants, despite their low  
360   frequency, contribute significantly to the genetic architecture of CAD and can help  
361   explain a portion of its missing heritability that common variants alone cannot account  
362   for.

363           There are several limitations in the study. First, there was a difference in age  
364   distribution between cases and controls. This discrepancy arose because we specifically  
365   selected early-onset CAD patients for the case group, resulting in a younger average age.  
366   As in previous rare variant studies, we prioritized selecting early-onset CAD cases to  
367   enrich genetic contributions <sup>20</sup>. Second, some of the prioritized genes for CAD in this

368 study have unknown functions, especially in cluster 6. However, many loci and genes  
369 identified in GWAS on CAD remain functionally uncharacterized, as well <sup>41,53</sup>.  
370 Therefore, future research is necessary to investigate the gene function and biological  
371 pathways to CAD development. Third, this study used WGS data from the Japanese  
372 population, so it is not certain whether the RVS created in this study can be applied to  
373 other populations since a PRS derived from GWAS in one population is reported to be  
374 less accurate in other populations <sup>11,54</sup>. These results need to be validated in other  
375 populations and prospective cohorts.

376           Taken together, our study underscores the important role of rare variants in the  
377 genetic landscape of CAD. By leveraging a machine learning-based framework, we  
378 have revealed CAD-associated genes and pathways influenced by rare variants. Our  
379 results demonstrate the distinct and complementary value of RVS compared to  
380 conventional PRS, highlighting the enhanced predictive power achieved through their  
381 integration. This comprehensive approach offers new insights into the pathogenesis of  
382 CAD, potentially leading to the accurate assessment and management of individual  
383 CAD risk.

384

385 **Consortia**

386 **The Biobank Japan Project**

387 Koichi Matsuda<sup>1,2</sup>, Takayuki Morisaki<sup>2,3</sup>, Yukinori Okada<sup>4</sup>, Yoichiro Kamatani<sup>5</sup>, Kaori  
388 Muto<sup>6</sup>, Akiko Nagai<sup>6</sup>, Yoji Sagiya<sup>2</sup>, Natsuhiko Kumasaka<sup>7</sup>, Yoichi Furukawa<sup>8</sup>, Yuji  
389 Yamanashi<sup>3</sup>, Yoshinori Murakami<sup>3</sup>, Yusuke Nakamura<sup>3</sup>, Wataru Obara<sup>9</sup>, Ken Yamaji<sup>10</sup>,  
390 Kazuhisa Takahashi<sup>11</sup>, Satoshi Asai<sup>12,13</sup>, Yasuo Takahashi<sup>13</sup>, Shinichi Higashiue<sup>14</sup>, Shuzo  
391 Kobayashi<sup>14</sup>, Hiroki Yamaguchi<sup>15</sup>, Yasunobu Nagata<sup>15</sup>, Satoshi Wakita<sup>15</sup>, Chikako Nito<sup>16</sup>,  
392 Yu-ki Iwasaki<sup>17</sup>, Shigeo Murayama<sup>18</sup>, Kozo Yoshimori<sup>19</sup>, Yoshio Miki<sup>20</sup>, Daisuke  
393 Obata<sup>21</sup>, Masahiko Higashiyama<sup>22</sup>, Akihide Masumoto<sup>23</sup>, Yoshinobu Koga<sup>23</sup> & Yukihiro  
394 Koretsune<sup>24</sup>

395

396 <sup>1</sup>Laboratory of Genome Technology, Human Genome Center, Institute of Medical  
397 Science, The University of Tokyo, Tokyo, Japan.

398 <sup>2</sup> Laboratory of Clinical Genome Sequencing, Graduate School of Frontier Sciences,  
399 The University of Tokyo, Tokyo, Japan.

400 <sup>3</sup>The Institute of Medical Science, The University of Tokyo, Tokyo, Japan.

401 <sup>4</sup>Department of Genome Informatics, Graduate School of Medicine, The University of  
402 Tokyo, Tokyo, Japan.

403 <sup>5</sup> Laboratory of Complex Trait Genomics, Graduate School of Frontier Sciences, The  
404 University of Tokyo, Tokyo, Japan.

405 <sup>6</sup> Department of Public Policy, Institute of Medical Science, The University of Tokyo,  
406 Tokyo, Japan.

407 <sup>7</sup> Division of Digital Genomics, Institute of Medical Science, The University of Tokyo,  
408 Tokyo, Japan.

409 <sup>8</sup> Division of Clinical Genome Research, Institute of Medical Science, The University of  
410 Tokyo, Tokyo, Japan.

411 <sup>9</sup> Department of Urology, Iwate Medical University, Iwate, Japan.

412 <sup>10</sup> Department of Internal Medicine and Rheumatology, Juntendo University Graduate  
413 School of Medicine, Tokyo, Japan.

414 <sup>11</sup> Department of Respiratory Medicine, Juntendo University Graduate School of  
415 Medicine, Tokyo, Japan.

416 <sup>12</sup> Division of Pharmacology, Department of Biomedical Science, Nihon University  
417 School of Medicine, Tokyo, Japan.

418 <sup>13</sup> Division of Genomic Epidemiology and Clinical Trials, Clinical Trials Research  
419 Center, Nihon University. School of Medicine, Tokyo, Japan.

420 <sup>14</sup> Tokushukai Group, Tokyo, Japan.

421 <sup>15</sup> Department of Hematology, Nippon Medical School, Tokyo, Japan.

422 <sup>16</sup> Laboratory for Clinical Research, Collaborative Research Center, Nippon Medical  
423 School, Tokyo, Japan.

424 <sup>17</sup> Department of Cardiovascular Medicine, Nippon Medical School, Tokyo, Japan.

425 <sup>18</sup> Tokyo Metropolitan Geriatric Hospital and Institute of Gerontology, Tokyo, Japan.

426 <sup>19</sup> Fukujuji Hospital, Japan Anti-Tuberculosis Association, Tokyo, Japan.

427 <sup>20</sup> The Cancer Institute Hospital of the Japanese Foundation for Cancer Research, Tokyo,  
428 Japan.

429 <sup>21</sup> Center for Clinical Research and Advanced Medicine, Shiga University of Medical  
430 Science, Shiga, Japan.

431 <sup>22</sup> Department of General Thoracic Surgery, Osaka International Cancer Institute, Osaka,  
432 Japan.

433 <sup>23</sup> Iizuka Hospital, Fukuoka, Japan.

434 <sup>24</sup> National Hospital Organization Osaka National Hospital, Osaka, Japan.

435

#### 436 **Acknowledgements**

437           We thank the staff of BBJ and the Nagahama cohort study for their assistance  
438 in collecting samples and clinical information. We thank the participants in the BBJ and  
439 Nagahama cohort study for their contribution to the study. H.I. is funded by the Japan  
440 Society for the Promotion of Science grant (JP22J00780, JP22K16128). K.I. is  
441 supported by the Japan Agency for Medical Research and Development (AMED) under  
442 grant numbers JP24bm1423005, JP24km0405209, JP24tm0524004, JP24tm0624002,  
443 JP24km0405209 and JP24ek0210164. K.I. and K.O. are supported by the Research  
444 Funding for Longevity Sciences from the NCGG (24–15). BBJ is supported by the  
445 Tailor-Made Medical Treatment Program of the Ministry of Education, Culture, Sports,  
446 Science, and Technology (MEXT) and AMED under grant numbers JP17km0305002  
447 and JP17km0305001, JP.24tm0624002. The Nagahama study was supported by a JSPS  
448 Grant-in-Aid for Scientific Research (C), KAKENHI grant numbers JP17K07255 and  
449 JP17KT0125, and the Practical Research Project for Rare/Intractable Diseases from  
450 AMED under grant numbers JP16ek0109070, JP18kk0205008, JP18kk0205001,  
451 JP19ek0109283, and JP19ek0109348.

452

#### 453 **Author contributions**

454           H.I. and K.I. conceived and designed the study. C.K., J.S., K.H., and F.M.  
455 collected, managed and genotyped the Nagahama cohort. K.M., C.T. and Y.K. collected  
456 and managed the BBJ samples. H.I. and K.I. analyzed WGS data, developed the

457 machine-learning model and performed the statistical analyses. S. K estimated the effect  
458 size to calculate PRS. S.Z. developed the PPI network module and analyzed it. H.Y.,  
459 R.K., H.M., K.M., N.E., K.O., Y.O., C.T., and Y.K. contributed to data processing,  
460 analysis and interpretation. K.I., M.S. and I.K. supervised the study. H.I. and K.I. wrote  
461 the manuscript, and many authors have provided valuable insights and edits.

462

### 463 **Declaration of Interests**

464 H.I. reports receiving grants from the Japan Heart Foundation / Bayer  
465 Pharmaceutical Research Grant Abroad. M.S. is a co-founder and the scientific advisory  
466 board member of Personalis, Qbio, January, SensOmics, Filtricine, Akna, Protos, Mirvie,  
467 NiMo, Onza, Oralome, Marble Therapeutics, and Iollo. He is also on the scientific  
468 advisory board of Danaher, Genapsys and Jupiter.

469

### 470 **Supplemental information**

471 **Document S1.** Table S1-S3, S5, S7, S9-13, Figure S1-8

472 **Table S4.** Summary statistics of aggregated rare variant association analysis using  
473 SAIGE-GENE+, related to Figure S4.

474 **Table S6.** Summary of 59 HEAL<sub>CAD</sub> genes with gene-based annotation, related to  
475 Figure 2.

476 **Table S8.** The 46 Protein Interaction Modules Identified in CAD, Related to Figure 2.

477

478

479 **Figure legends**

480 **Figure 1. Overview of the current study.**

481 We studied the genetic factors of coronary artery disease (CAD) combining whole-  
482 genome sequencing data and a machine learning-based framework named the modified  
483 HEAL method in patients with MI, one of the most severe forms of CAD, and controls.  
484 We sequenced the whole genomes of Japanese CAD patients and controls and applied  
485 the modified HEAL method framework. The framework was based on a sparse  
486 modeling devised to distinguish diseased individuals from controls. After the  
487 hyperparameter tuning and training of the model by the cross-validation method, the  
488 model outputted a list of genes related to CAD, which were subsequently analyzed by a  
489 clustering-based method and mapped on the protein-protein interaction network to  
490 reveal the CAD-associated modules. The function of the identified genes was also  
491 confirmed by the human phenotype and knockout mouse phenotype databases. The  
492 learned (optimized) machine learning model was applied to derive rare variant-based  
493 genetic risk scores (RVS) to predict CAD outcomes in an independent validation cohort.  
494 We also tested the relationship of the RVS with clinical features and common variant-  
495 based polygenic risk score (PRS). RVS was combined with PRS to improve the  
496 prediction performance of CAD disease status in the independent validation cohort. BBJ,  
497 BioBank Japan; MI, myocardial infarction; CRS, combined risk score

498

499

500 **Figure 2. Functional analysis of HEAL<sub>CAD</sub> genes**

501 (A) Fifty-nine genes identified by the machine learning-based framework were  
502 annotated using six different criteria; 1) The constraint score (pLI) from the gnomAD  
503 database 2) Overlap with GWAS on CAD and its risk factor (lipids, diabetes, obesity,  
504 blood pressure, coagulation, smoking) phenotypes, 3) Overlap with the genes in which  
505 GWAS-significant variants act as enhancers, 4) Knock-out mouse phenotype with blood  
506 pressure, diabetes, and lipid traits, 5) Human phenotype ontology and 6) Gene ontology.  
507 Then the fifty-nine genes were grouped into eight clusters by hierarchical clustering  
508 based on functional annotations. For GWAS and Genehancer, red indicates a significant  
509 association and light red denotes suggestive significance. (B) Gene ontology (GO) and  
510 human phenotype ontology (HPO) term enrichment analysis. The GO and HPO  
511 annotation results were based on 59 genes. Gene ontology categories included  
512 molecular function, cellular components and biological process. GO and HPO  
513 categories for each function were sorted by decreasing order of evidence based on the  
514 GO enrichment test P-value. Only the significant categories after multiple test  
515 corrections are shown. (C) The forty-six modules were identified in the protein-protein-  
516 interaction network using diffusion component analysis seeded by the 59 HEAL<sub>CAD</sub>  
517 genes. (D) Visualization of the module 119 network of the protein-protein interactions.  
518 The module included important genes involved in cholesterol metabolism, including  
519 *LDLR*, *PCSK9*, *ANGPTL3*, *ANGPTL4*, and *LIPA*. GWAS, genome-wide association  
520 study; CAD, coronary artery disease; DM, diabetes mellitus; BP, blood pressure; IMPC,  
521 International Mouse Phenotyping Consortium; HP, human phenotype; GOMF, gene  
522 ontology molecular function; GOBP, gene ontology biological pathway; GOCC, gene  
523 ontology cellular component.



524

525 **Figure 3. Rare variant risk score (RVS) and its clinical impact**

526 (A) Correlation between RVS and continuous clinical indices. Data are presented as  
527 Pearson's correlation coefficients and their 95% confidence intervals (CIs). Exact P  
528 values are shown in Table S9. (B) Kaplan-Meier curves for cardiovascular mortality  
529 among total participants stratified into two groups based on RVS. Participants with high  
530 RVS died significantly earlier than those with low RVS. (C) Kaplan-Meier curves for  
531 cardiovascular mortality among CAD patients (n=200) stratified into two groups based  
532 on RVS. CAD patients with high RVS (top 5%) showed significantly worse  
533 cardiovascular prognosis. LDLC, low-density lipoprotein cholesterol; Tbil, total  
534 bilirubin; ALT, alanine aminotransferase; PTINR, prothrombin time international  
535 normalized ratio; TC, total cholesterol; K, potassium; Hb, hemoglobin; UA, uric acid;  
536 APTT, activated partial thromboplastin time; Alb, albumin; RBC, red blood cell; AST,  
537 aspartate aminotransferase; WBC, white blood cell; CK, creatine kinase; TP, total  
538 protein; Cre, creatinine; DBP, diastolic blood pressure; SBP, systolic blood pressure;  
539 BUN, blood urea nitrogen; TG, triglycerides; CRP, C-reactive protein; PLT, platelet; P,  
540 Phosphorus;  $\gamma$ GTP, gamma-glutamyl transpeptidase; BS, blood sugar; LDH, Lactate  
541 dehydrogenase.

542

543

544 **Figure 4. The Relationship between RVS, PRS, CRS, and clinical indices.**

545 (A) A scatter plot illustrating the relationship between RVS and PRS, with cases (red)  
546 and controls (gray) color-coded. The overall (gray) and case-only (pink) regression lines  
547 and correlation coefficients are shown. A significant negative correlation was observed  
548 in the CAD cases. (B) Correlation between combined risk score (CRS), defined by the  
549 average of RVS and PRS, and continuous clinical indices. Data are presented as  
550 Pearson's correlation coefficients and their 95% CIs. Exact P values are shown in Table  
551 S11. (C) Correlation between clinical measurements and different genetic risk scores  
552 (RVS, PRS and CRS). Only significant correlations are displayed with a circle. Blue,  
553 positive correlation; red, negative correlation. Larger circles correspond to a stronger  
554 correlation. LDLC, low density lipoprotein cholesterol; Tbil, total bilirubin; ALT,  
555 alanine aminotransferase; PTINR, prothrombin time international normalized ratio; TC,  
556 total cholesterol; K, potassium; Hb, hemoglobin; UA, uric acid; APTT, activated partial  
557 thromboplastin time; Alb, albumin; RBC, red blood cell; AST, aspartate  
558 aminotransferase; WBC, white blood cell; CK, creatine kinase; TP, total protein; Cre,  
559 creatinine; DBP, diastolic blood pressure; SBP, systolic blood pressure; BUN, blood  
560 urea nitrogen; TG, triglycerides; CRP, C-reactive protein; PLT, platelet; P, Phosphorus;  
561  $\gamma$ GTP, gamma-glutamyl transpeptidase; BS, blood sugar; LDH, Lactate dehydrogenase  
562

563 **Figure 5. The combined RVS and PRS risk score improved CAD prediction**

564 (A) Receiver operating characteristic (ROC) curve for RVS, PRS and CRS (Combined  
565 Risk Score). The curve plots the true positive rate (sensitivity) against the false positive  
566 rate (1-specificity) for different threshold values of the predictive score. The area under  
567 the curve (AUC) is indicated, representing the score's accuracy in predicting the  
568 outcome. The dotted line represents a reference line of no discrimination (AUC = 0.5).  
569 Points on the curve closer to the top-left corner indicate higher diagnostic accuracy. (B)  
570 Precision-recall curve (PRC) for RVS, PRS and CRS. The curve shows the trade-off  
571 between precision (positive predictive value) and recall (sensitivity) at various threshold  
572 levels. The confidence interval for the area under the PRC was estimated from the  
573 20,000 times bootstrap replication method. (C) Boxplot of Pseudo  $R^2$  for CAD  
574 prediction performance. This box plot displays the pseudo- $R^2$  values comparing the  
575 CAD prediction performance of RVS, PRS and CRS. The distribution of pseudo- $R^2$  was  
576 estimated from 20,000 times bootstrapping. The box plot center line represents the  
577 median, the bounds represent the first and third quartile, and the whiskers reach to 1.5  
578 times the interquartile range.

579 **Tables**

580 **Table 1. Demographic features of participants**

Data	Disease status	Total N	Males		Age (years)		BMI (kg/m <sup>2</sup> )		Age at MI onset (years)	
			N	%	Mean	SD	Mean	SD	Mean	SD
Discovery cohort	Case	1,752	1,617	92.29	60.1	13.7	25.0	3.4	47.4	4.1
	Control	3,019	1,205	39.91	55.3	8.0	23.9	4.0	-	-
Validation cohort	Case	200	183	91.50	43.8	9.8	26.5	4.3	36.0	3.9
	Control	824	420	50.97	49.3	13.0	22.9	3.8	-	-

581 **SD, standard deviation; BMI, body mass index; MI, myocardial infarction**

582

583 **STAR Methods**

584 **Code availability**

585 The code of the modified HEAL framework is available on  
586 <https://github.com/pirocv/HEAL>.

587 **Study cohort**

588 Two previously described cohorts were used in the current study. BioBank  
589 Japan (BBJ) is a hospital-based Japanese biobank project including clinical and genetic  
590 data from a variety of patients<sup>55,56</sup>. Participants were recruited from 12 hospitals  
591 throughout Japan. The Nagahama Prospective Genome Cohort (Nagahama study) is the  
592 genome cohort conducted in Shiga, Japan. Participants aged 30–74 years were recruited  
593 from the general population in Nagahama city from 2007 to 2010<sup>57</sup>.

594

595 **Whole genome sequencing and quality control**

596 We sequenced 1,765 CAD patients and 3,148 controls from the cohort. Whole  
597 genome sequence (WGS) was performed on Illumina’s HiSeqX aiming at 15x depth,  
598 using 150-base pair-end reads. We also sequenced an additional 200 CAD cases and 836  
599 controls aiming at 30x depth using 150-base paired-end reads. In order to enrich for a  
600 genetic contribution to disease<sup>20</sup>, we prioritized patients with early-onset MI, one of the  
601 most severe forms of CAD, within the BBJ cohort for WGS (age of MI onset in 15x and  
602 30x WGS cohort:  $47.4 \pm 4.1$  years and  $36.0 \pm 3.9$  years, respectively). Sequenced reads  
603 were aligned to the hs37d5 reference genome using BWA software<sup>58</sup>. The genotypes of  
604 the samples were called using the HaplotypeCaller implemented in GATK v3.8. Per-  
605 sample Genomic Variant Call Format (gVCF) genotype data were merged and jointly  
606 called using GenotypeGVCFs. We defined exclusion filters for genotypes as follows.

607 (1) For 15x depth data, filtered depth (DP) < 2, quality of the assigned genotype  
608 (genotype quality; GQ) < 20. (2) For 30x depth data, DP < 5, GQ < 20, DP > 60 and GQ  
609 < 95. We set these genotypes as missing and excluded variants with call rates < 90%  
610 before variant quality score recalibration. For sample quality control, the following  
611 samples were excluded: (1) age < 20 years old, (2) excess missing genotypes (> 10%),  
612 (3) samples whose genetically inferred sex did not match the self-reported sex, (4)  
613 closely related samples estimated by identity-by-descent and identity-by-state analysis  
614 ( $\hat{\pi} > 0.1875$ ) and (5) excess heterozygosity. We also excluded non-Japanese  
615 participants estimated from Principal component analysis (PCA) calculated using  
616 PLINK 2.0<sup>59</sup>. The total number of genomes that failed data quality control is  
617 summarized in **Table S13**. After the sample quality control, we retained 1,752 CAD  
618 case samples and 3019 non-CAD control samples for 15x depth data and 200 case  
619 samples and 824 control samples for 30x depth. Then, the variant quality control was  
620 performed excluding (1) high missingness (5% for 15x depth and 1% for 30x depth), (2)  
621 Hardy-Weinberg equilibrium ( $P < 1 * 10^{-6}$ ), (3) variants in the low complexity region.  
622 WGS data with 15x depth data was used as a discovery cohort and the 30x depth data  
623 was used as the validation cohort in the machine learning-based analysis.

624

### 625 **Single variant association analysis**

626 The single variant association test was performed by logistic regression  
627 implemented in PLINK 2.0<sup>60</sup> with adjustment for age, sex, and the first 10 principal  
628 components of ancestry. Principal components of ancestry were calculated using PLINK  
629 2.0<sup>59</sup>. The inclusion of principal components as covariates in the logistic regression  
630 analysis increases the power to detect true genetic associations and minimizes

631 confounding by population stratification<sup>61</sup>. Variants with a missing rate of less than  
632 0.01 were included in the analysis. Genomic inflation factor ( $\lambda_{GC}$ ) was calculated using  
633 variants with  $MAF \geq 0.001$ . Single variant association analysis was also performed  
634 using SAIGE<sup>62</sup> with adjustment for age, sex, and the first 10 principal components of  
635 ancestry. SAIGE is widely used in GWASs for binary traits to account for population  
636 structure and relatedness while correcting for the type I error rates<sup>62</sup>. The genome-wide  
637 significance threshold was set at  $P = 5 * 10^{-8}$ . To define a locus, we added 500 kb to  
638 both sides of each genome-wide significant SNP and merged overlapping regions. To  
639 determine whether each locus was novel, a literature search was conducted to ascertain  
640 if any of the regions contained SNPs had been previously reported as significant for  
641 CAD.

642

#### 643 **Aggregated rare variant association analysis**

644 We also performed gene-based association analysis using SAIGE-GENE+  
645 software, which accounts for the relatedness among the study samples<sup>63,64</sup>. We first  
646 calculated sparse GRM using the WGS data and fit the null model in the SAIGE-  
647 GENE+ algorithm step1. For the gene-based association analysis, we extracted rare  
648 ( $MAF < 0.001$ ) nonsynonymous variants including (nonsynonymous single nucleotide  
649 variations (SNV), nonframeshift insertion, nonframeshift deletion, frameshift insertion,  
650 frameshift deletion, stopgain, stoploss, and splice site variants). Splice-site variants,  
651 pLOF variants and damaging missense variants defined by a REVEL score  $> 0.5$ <sup>65</sup> were  
652 included in the analysis. SKAT-O test implemented in SAIGE-GENE+ software was  
653 performed with adjustment for age, sex and first 10 principal components of ancestry.  
654 Gene-wide significance threshold and suggestive threshold were set at  $P = 2.5 * 10^{-6}$



655 and  $P = 5 * 10^{-4}$ , respectively. Statistical inflation was estimated by Q-Q plot.

656

### 657 **Machine learning-based analysis (modified HEAL)**

658 We employed a recently developed machine learning-based rare variant  
659 analysis method called HEAL (hierarchical estimate from agnostic learning). A detailed  
660 HEAL method is described in the original paper<sup>19</sup>. In this framework (**Figure S8**), we  
661 first annotated each variant using ANNOVAR software<sup>66</sup> and extracted rare  
662 nonsynonymous variants (nonsynonymous SNV, nonframeshift insertion, nonframeshift  
663 deletion, frameshift insertion, frameshift deletion, stopgain, stoploss, and splice site  
664 variants) that were not present in the East-Asian populations analyzed in the 1000  
665 Genomes Project<sup>67</sup>. Variants with high frequency in the WGS data and gnomAD East  
666 Asian database<sup>68</sup> ( $MAF > 0.1$ ) were also filtered. To estimate the mutation burden for  
667 each gene based on the rare variants, we used the REVEL score (ranges from 0 to 1 with  
668 a higher score indicating a damaging variant), which was internally computed by  
669 ANNOVAR software. The deleteriousness score of the putative loss of function (pLOF)  
670 variants, such as stopgain and splice site variants, was set as 1. Next, we calculated the  
671 cumulative effects of rare nonsynonymous variants for each gene as

$$g_{in} = \sum_{j=1}^{m_{in}} s_{ijn}$$

672 , where  $g_{in}$  is the mutation burden of the gene  $i$  of  $n$ th sample,  $m_{in}$  is the number of rare  
673 nonsynonymous variants,  $s_{ijn}$  is the deleteriousness score for variant  $j$  of gene  $i$ . Using  
674 the above formula, we obtained a matrix of estimated mutation burden for each gene per  
675 sample ( $\mathbf{x}_n = (g_{1n}, g_{2n}, \dots, g_{mn})$ , where  $m$  is the number of the total genes). The  
676 mutation burden was standardized (Z-score normalization). We trained a regularized

677 logistic regression model for a genome-based CAD prediction model. The input of the  
678 model is the calculated mutation burden and the output is the probability of CAD as  
679 shown in the following equation.

$$\widehat{y}_n = P(y_n = 1 | \mathbf{x}_n) = \sigma(\mathbf{w}^T \mathbf{x}_n) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}_n)}$$

680 , where  $y_n$  is the label for CAD case (1) or control (0),  $\widehat{y}_n$  is the probability of being  
681 CAD positive given the mutation burden  $\mathbf{x}_n$  for the  $n$ th sample,  $\sigma$  is the sigmoid  
682 function and  $\mathbf{w}$  is the weight vector. To identify the optimal coefficient vector  $\mathbf{w}$  that  
683 achieve the maximum consistency between the model probabilities ( $\widehat{y}_n$ ) and the  
684 observations for the cohort ( $y_n$ ), we solved the following optimization problem.

$$\min_{\mathbf{w}} -\frac{1}{N} \sum_{n=1}^N y_n \log \widehat{y}_n + (1 - y_n) \log(1 - \widehat{y}_n) + \lambda \|\mathbf{w}\|_1$$

685 In this regularized logistic regression, regularization strength is determined by  
686 parameter  $\lambda$ , and it is a hyperparameter of the machine learning model, which was  
687 determined by the cross-validation method (**Figure S5**). By training the model to  
688 predict disease status, it outputs the minimal set of most distinguishing features (genes)  
689 for CAD. The trained model can be used to estimate the rare variant-based disease risk  
690 score (RVS) from the genomic data. We have named this the modified HEAL because  
691 our approach differs from the original method in that we included not only missense  
692 variants but also pLOF variants. We determined the hyperparameters using grid search  
693 and estimated the performance in the independent cohort to avoid bias and  
694 overestimation of the model's performance, while the performance was estimated using  
695 internal cross-validation in the original method.

696

## 697 **Interpretation of genes identified by modified HEAL**

698 To investigate the functions of the 59 identified genes, we first annotated each  
699 one using various databases and then conducted clustering analysis to categorize the  
700 groups of genes to obtain the eight functional groups. Annotations included checking  
701 the constraint score (pLI) from the gnomAD database<sup>68</sup>, identifying whether the genes  
702 were reported in previous GWAS on CAD and its risk factors (lipids, diabetes, obesity,  
703 blood pressure, coagulation function, and smoking-related phenotypes) using the  
704 GWAS Catalog, and checking for the overlap with target genes of enhancers that were  
705 significant in previous GWAS on CAD and its risk factors (same as above) using the  
706 GeneHancer database, which includes genome-wide enhancers and their target genes<sup>21</sup>.  
707 Further analysis involved examining the International Mouse Phenotyping Consortium  
708 (IMPC) database to determine if the corresponding genes in knock-out mice are  
709 significantly related to phenotypes such as blood pressure, blood glucose and lipid traits.  
710 Enrichment analysis for Gene Ontology and Human Phenotype Ontology was  
711 performed using g:Profiler<sup>69</sup> to gain insights into the biological processes and human  
712 phenotypic abnormalities associated with these genes<sup>22</sup>. We considered statistical  
713 significance for the enrichment analysis with a false discovery rate under 0.1.

714 To analyze the functional modules in CAD, we downloaded the human  
715 protein-protein interactions (PPIs) from STRING v12.0, comprising 19,622 proteins and  
716 6,857,702 interactions. High-confidence PPIs (combined score >700) were extracted for  
717 downstream analysis, including 16,185 proteins and 236,000 interactions. To remove  
718 bias from hub proteins, we applied the random walk with restart (RWR) algorithm with  
719 a restart probability of 0.5. This produced a smoothed network after retaining the top  
720 5% predicted edges (n = 6,243,766). We employed the Louvain method<sup>70</sup> to decompose

721 the network into different modules. Following algorithm convergence, we obtained  
722 1,261 modules with an average size of 13 nodes. Among the 1,261 PPI modules, 46  
723 encompassed at least one gene identified by the machine learning analysis. We used  
724 g:profiler to determine functional enrichment for each module. Cytoscape software <sup>71</sup>  
725 was used to visualize the PPI modules.

726

### 727 **Genetic risk scores**

728 By optimizing the machine learning-based model, the modified HEAL  
729 framework can also make a prediction of disease based on the input genome. We call it  
730 rare variant-based genetic risk score (RVS) because it only leverages information on  
731 rare variants. Using the trained model, we estimated the RVS prediction performance in  
732 the validation cohort. We also analyzed the association between RVS and clinical  
733 parameters such as vital signs and blood test data in the BBJ data using Pearson's  
734 correlation. To investigate the prognostic impact of RVS, we divided the patients into  
735 those in the top 5% and those below, then compared their outcome using Kaplan-Meier  
736 analysis and a log-rank test. To compare the properties between RVS and the common  
737 variant-based polygenic risk score (PRS), GWAS of CAD in BBJ (case 25,668 vs  
738 control 141,667) was performed. The individuals included in the GWAS were  
739 genotyped using the HumanOmniExpressExome v.1.0/v.1.2 platform (Illumina) or in  
740 combination with HumanOmniExpress v.1.0 and Human Exome BeadChip v.1.0/v.1.1  
741 (Illumina). For genotype quality control, variants with (1) SNP call rate < 99%, (2)  
742 Hardy-Weinberg equilibrium ( $P < 1 * 10^{-6}$ ) and (3) heterozygous counts <5 were  
743 excluded. We performed pre-phasing using Eagle software. Phased haplotypes were  
744 imputed to the in-house reference panel from BBJ <sup>11</sup> by minimac3 <sup>72</sup>. Variants with low

745 imputation quality ( $R^2 < 0.3$ ) were excluded. GWAS was performed by logistic  
746 regression implemented in PLINK 2.0<sup>60</sup> with adjustment for age, age<sup>2</sup>, sex and first 10  
747 principal components of ancestry. Then PRS of *i*th sample was calculated as follows

$$PRS_i = \sum_{j=1}^M a_{i,j} \beta_j$$

748 , where *M* is the number of variants in GWAS,  $a_{i,j}$  is the number of effect allele of *j*th  
749 variant in *i*th sample, and  $\beta_j$  is the effect size of *j*th variant estimated by GWAS. The  
750 number of variants included in the PRS calculation was determined by the pruning and  
751 thresholding method<sup>13</sup>. The relationship between RVS and PRS was examined by  
752 Pearson's correlation coefficient, both in cases only and across the validation cohort. We  
753 then integrated both RVS and PRS by normalizing (mean 0, standard deviation 1) and  
754 adding them together to obtain combined risk score (CRS). The predictive performance  
755 of each genetic score was estimated on the validation cohort, which was not used in the  
756 derivation of the RVS, PRS, or CRS. We used receiver operating characteristics (ROC)  
757 to evaluate the predictive performance. To examine whether CRS improves predictive  
758 performance compared to conventional PRS, we compared AUROC of PRS and CRS  
759 by DeLong's test. We also calculated the area under precision-recall curve (AUPRC)  
760 and Nagelkerke's pseudo  $R^2$  metrics. The *P* values were derived using a 20000 times  
761 bootstrap replication method. In all statistical analyses, R software was used and a two-  
762 sided  $P < 0.05$  was considered statistically significant.

763

764

## 765 Reference

- 766 1. Wang, H., Naghavi, M., Allen, C., Barber, R.M., Bhutta, Z.A., Carter, A., Casey,  
767 D.C., Charlson, F.J., Chen, A.Z., Coates, M.M., et al. (2016). Global, regional, and  
768 national life expectancy, all-cause mortality, and cause-specific mortality for 249  
769 causes of death, 1980–2015: a systematic analysis for the Global Burden of  
770 Disease Study 2015. *Lancet* 388, 1459–1544.
- 771 2. Roth, G.A., Huffman, M.D., Moran, A.E., Feigin, V., Mensah, G.A., Naghavi, M.,  
772 and Murray, C.J.L. (2015). Global and regional patterns in cardiovascular mortality  
773 from 1990 to 2013. *Circulation* 132, 1667–1678.
- 774 3. McPherson, R., and Tybjaerg-Hansen, A. (2016). Genetics of Coronary Artery  
775 Disease. *Circ. Res.* 118, 564–578.
- 776 4. Musunuru, K., and Kathiresan, S. (2019). Genetics of Common, Complex  
777 Coronary Artery Disease. *Cell* 177, 132–145.
- 778 5. Khera, A.V., and Kathiresan, S. (2017). Genetics of coronary artery disease:  
779 discovery, biology and clinical translation. *Nat. Rev. Genet.* 18, 331–344.
- 780 6. Marenberg, M.E., Risch, N., Berkman, L.F., Floderus, B., and de Faire, U. (1994).  
781 Genetic Susceptibility to Death from Coronary Heart Disease in a Study of Twins.  
782 *N. Engl. J. Med.* 330, 1041–1046.
- 783 7. Zdravkovic, S., Wienke, A., Pedersen, N.L., Marenberg, M.E., Yashin, A.I., and De  
784 Faire, U. (2002). Heritability of death from coronary heart disease: a 36-year  
785 follow-up of 20 966 Swedish twins. *J. Intern. Med.* 252, 247–254.
- 786 8. Nikpay, M., Goel, A., Won, H.-H., Hall, L.M., Willenborg, C., Kanoni, S.,  
787 Saleheen, D., Kyriakou, T., Nelson, C.P., Hopewell, J.C., et al. (2015). A  
788 comprehensive 1,000 Genomes-based genome-wide association meta-analysis of  
789 coronary artery disease. *Nat. Genet.* 47, 1121–1130.
- 790 9. van der Harst, P., and Verweij, N. (2018). Identification of 64 Novel Genetic Loci  
791 Provides an Expanded View on the Genetic Architecture of Coronary Artery  
792 Disease. *Circ. Res.* 122, 433–443.
- 793 10. Matsunaga, H., Ito, K., Akiyama, M., Takahashi, A., Koyama, S., Nomura, S., Ieki,  
794 H., Ozaki, K., Onouchi, Y., Sakaue, S., et al. (2020). Transethnic Meta-Analysis of  
795 Genome-Wide Association Studies Identifies Three New Loci and Characterizes  
796 Population-Specific Differences for Coronary Artery Disease. *Circulation:  
797 Genomic and Precision Medicine* 13, e002670.
- 798 11. Koyama, S., Ito, K., Terao, C., Akiyama, M., Horikoshi, M., Momozawa, Y.,  
799 Matsunaga, H., Ieki, H., Ozaki, K., Onouchi, Y., et al. (2020). Population-specific  
800 and trans-ancestry genome-wide analyses identify distinct and shared genetic risk  
801 loci for coronary artery disease. *Nat. Genet.* 52, 1169–1177.

- 802 12. Schnitzler, G.R., Kang, H., Fang, S., Angom, R.S., Lee-Kim, V.S., Ma, X.R., Zhou,  
803 R., Zeng, T., Guo, K., Taylor, M.S., et al. (2024). Convergence of coronary artery  
804 disease genes onto endothelial cell programs. *Nature* 626, 799–807.
- 805 13. Khera, A.V., Chaffin, M., Aragam, K.G., Haas, M.E., Roselli, C., Choi, S.H.,  
806 Natarajan, P., Lander, E.S., Lubitz, S.A., Ellinor, P.T., et al. (2018). Genome-wide  
807 polygenic scores for common diseases identify individuals with risk equivalent to  
808 monogenic mutations. *Nat. Genet.* 50, 1219–1224.
- 809 14. Torkamani, A., Wineinger, N.E., and Topol, E.J. (2018). The personal and clinical  
810 utility of polygenic risk scores. *Nat. Rev. Genet.* 19, 581–590.
- 811 15. Richardson, T.G., Harrison, S., Hemani, G., and Davey Smith, G. (2019). An atlas  
812 of polygenic risk score associations to highlight putative causal relationships  
813 across the human phenome. *Elife* 8. 10.7554/eLife.43657.
- 814 16. Minikel, E.V., Karczewski, K.J., Martin, H.C., Cummings, B.B., Whiffin, N.,  
815 Rhodes, D., Alföldi, J., Trembath, R.C., van Heel, D.A., Daly, M.J., et al. (2020).  
816 Evaluating drug targets through human loss-of-function genetic variation. *Nature*  
817 581, 459–464.
- 818 17. Backman, J.D., Li, A.H., Marcketta, A., Sun, D., Mbatchou, J., Kessler, M.D.,  
819 Benner, C., Liu, D., Locke, A.E., Balasubramanian, S., et al. (2021). Exome  
820 sequencing and analysis of 454,787 UK Biobank participants. *Nature* 599, 628–  
821 634.
- 822 18. Jurgens, S.J., Choi, S.H., Morrill, V.N., Chaffin, M., Pirruccello, J.P., Halford, J.L.,  
823 Weng, L.-C., Nauffal, V., Roselli, C., Hall, A.W., et al. (2022). Analysis of rare  
824 genetic variation underlying cardiometabolic diseases and traits among 200,000  
825 individuals in the UK Biobank. *Nat. Genet.* 54, 240–250.
- 826 19. Li, J., Pan, C., Zhang, S., Spin, J.M., Deng, A., Leung, L.L.K., Dalman, R.L., Tsao,  
827 P.S., and Snyder, M. (2018). Decoding the Genomics of Abdominal Aortic  
828 Aneurysm. *Cell* 174, 1361-1372.e10.
- 829 20. Do, R., Stitzel, N.O., Won, H.-H., Jørgensen, A.B., Duga, S., Angelica Merlini, P.,  
830 Kiezun, A., Farrall, M., Goel, A., Zuk, O., et al. (2015). Exome sequencing  
831 identifies rare LDLR and APOA5 alleles conferring risk for myocardial infarction.  
832 *Nature* 518, 102–106.
- 833 21. Fishilevich, S., Nudel, R., Rappaport, N., Hadar, R., Plaschkes, I., Iny Stein, T.,  
834 Rosen, N., Kohn, A., Twik, M., Safran, M., et al. (2017). GeneHancer: genome-  
835 wide integration of enhancers and target genes in GeneCards. *Database* 2017.  
836 10.1093/database/bax028.
- 837 22. Groza, T., Gomez, F.L., Mashhadi, H.H., Muñoz-Fuentes, V., Gunes, O., Wilson,  
838 R., Cacheiro, P., Frost, A., KeskiVali-Bond, P., Vardal, B., et al. (2023). The  
839 International Mouse Phenotyping Consortium: comprehensive knockout

- 840 phenotyping underpinning the study of human disease. *Nucleic Acids Res.* *51*,  
841 D1038–D1045.
- 842 23. Iacocca, M.A., Chora, J.R., Carrié, A., Freiburger, T., Leigh, S.E., Defesche, J.C.,  
843 Kurtz, C.L., DiStefano, M.T., Santos, R.D., Humphries, S.E., et al. (2018). ClinVar  
844 database of global familial hypercholesterolemia-associated DNA variants. *Hum.*  
845 *Mutat.* *39*, 1631–1640.
- 846 24. Versmissen, J., Oosterveer, D.M., Yazdanpanah, M., Dehghan, A., Hólm, H.,  
847 Erdman, J., Aulchenko, Y.S., Thorleifsson, G., Schunkert, H., Huijgen, R., et al.  
848 (2015). Identifying genetic risk variants for coronary heart disease in familial  
849 hypercholesterolemia: an extreme genetics approach. *Eur. J. Hum. Genet.* *23*, 381–  
850 387.
- 851 25. Tajima, T., Morita, H., Ito, K., Yamazaki, T., Kubo, M., Komuro, I., and  
852 Momozawa, Y. (2018). Blood lipid-related low-frequency variants in LDLR and  
853 PCSK9 are associated with onset age and risk of myocardial infarction in Japanese.  
854 *Sci. Rep.* *8*, 1–9.
- 855 26. Dickinson, M.E., Flenniken, A.M., Ji, X., Teboul, L., Wong, M.D., White, J.K.,  
856 Meehan, T.F., Weninger, W.J., Westerberg, H., Adissu, H., et al. (2016). High-  
857 throughput discovery of novel developmental phenotypes. *Nature* *537*, 508–514.
- 858 27. Huang, J., Huffman, J.E., Huang, Y., Do Valle, Í., Assimes, T.L., Raghavan, S.,  
859 Voight, B.F., Liu, C., Barabási, A.-L., Huang, R.D.L., et al. (2022). Genomics and  
860 phenomics of body mass index reveals a complex disease network. *Nat. Commun.*  
861 *13*, 7973.
- 862 28. Zhu, Z., Guo, Y., Shi, H., Liu, C.-L., Panganiban, R.A., Chung, W., O'Connor, L.J.,  
863 Himes, B.E., Gazal, S., Hasegawa, K., et al. (2020). Shared genetic and  
864 experimental links between obesity-related traits and asthma subtypes in UK  
865 Biobank. *J. Allergy Clin. Immunol.* *145*, 537–549.
- 866 29. Liu, M., Jiang, Y., Wedow, R., Li, Y., Brazel, D.M., Chen, F., Datta, G., Davila-  
867 Velderrain, J., McGuire, D., Tian, C., et al. (2019). Association studies of up to 1.2  
868 million individuals yield new insights into the genetic etiology of tobacco and  
869 alcohol use. *Nat. Genet.* *51*, 237–244.
- 870 30. Cecil, J.E., Tavendale, R., Watt, P., Hetherington, M.M., and Palmer, C.N.A.  
871 (2008). An Obesity-Associated FTO Gene Variant and Increased Energy Intake in  
872 Children. *N. Engl. J. Med.* *359*, 2558–2566.
- 873 31. Claussnitzer, M., Dankel, S.N., Kim, K.-H., Quon, G., Meuleman, W., Haugen, C.,  
874 Glunk, V., Sousa, I.S., Beaudry, J.L., Puviondran, V., et al. (2015). FTO Obesity  
875 Variant Circuitry and Adipocyte Browning in Humans. *N. Engl. J. Med.* *373*, 895–  
876 907.
- 877 32. Locke, A.E., Kahali, B., Berndt, S.I., Justice, A.E., Pers, T.H., Day, F.R., Powell, C.,

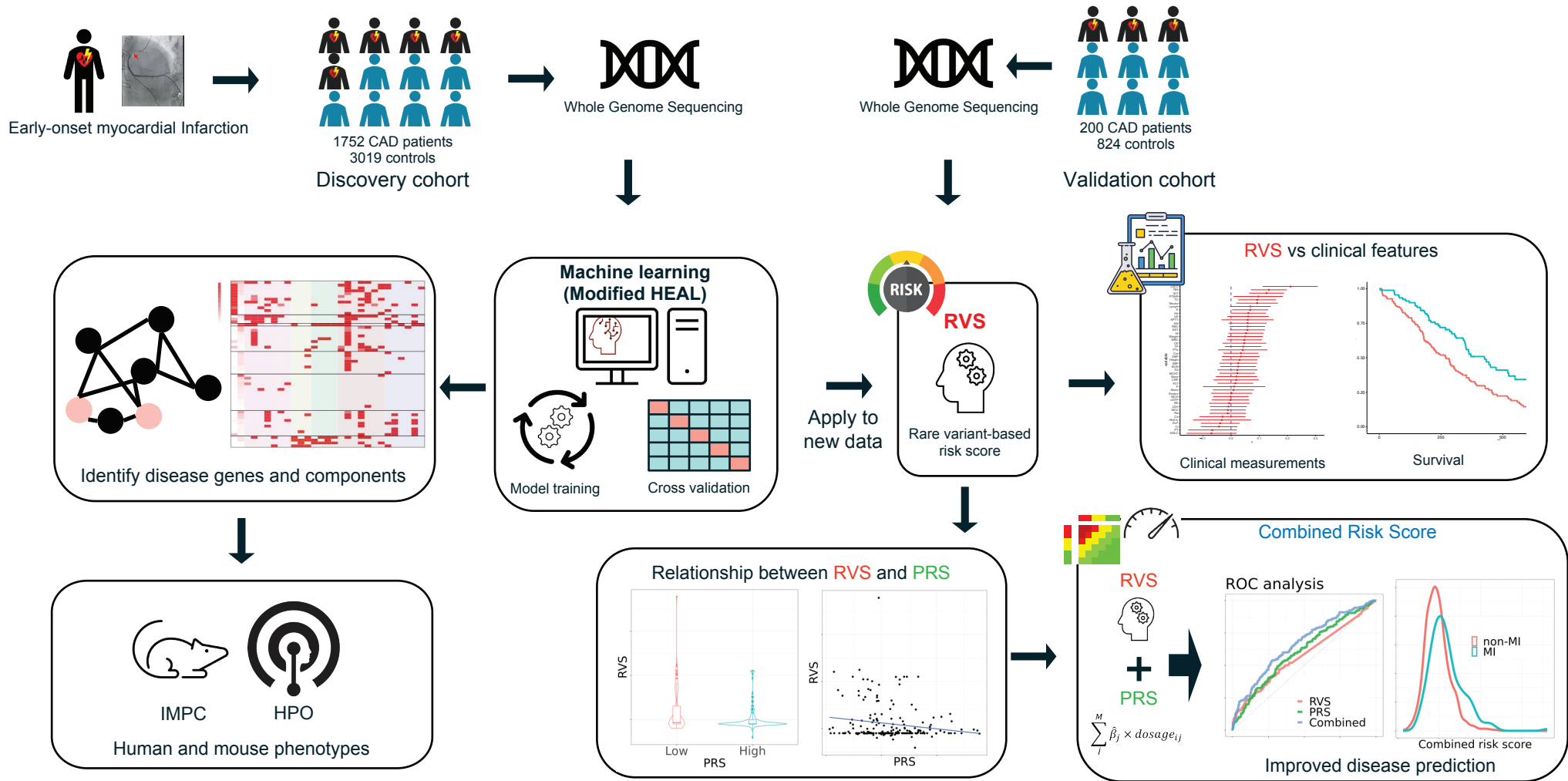


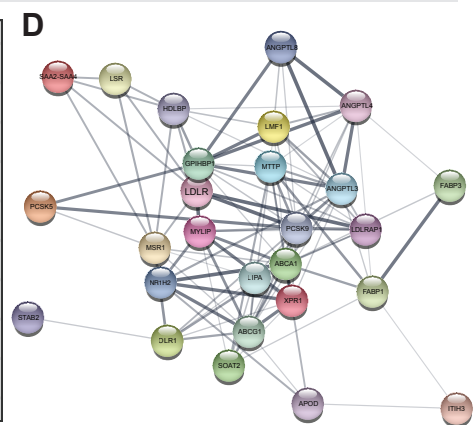
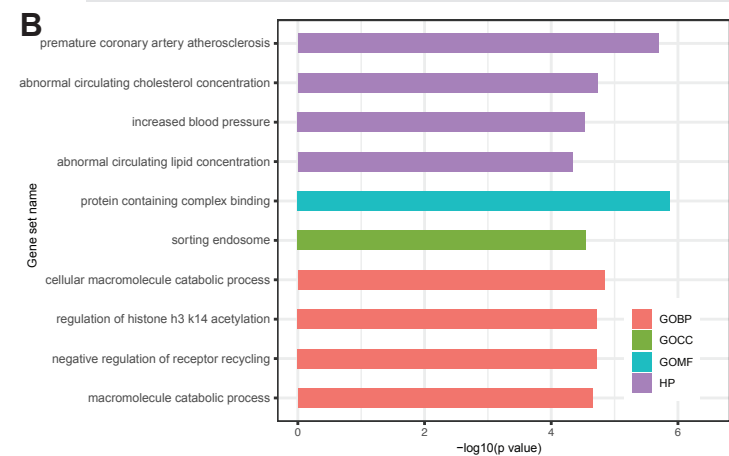
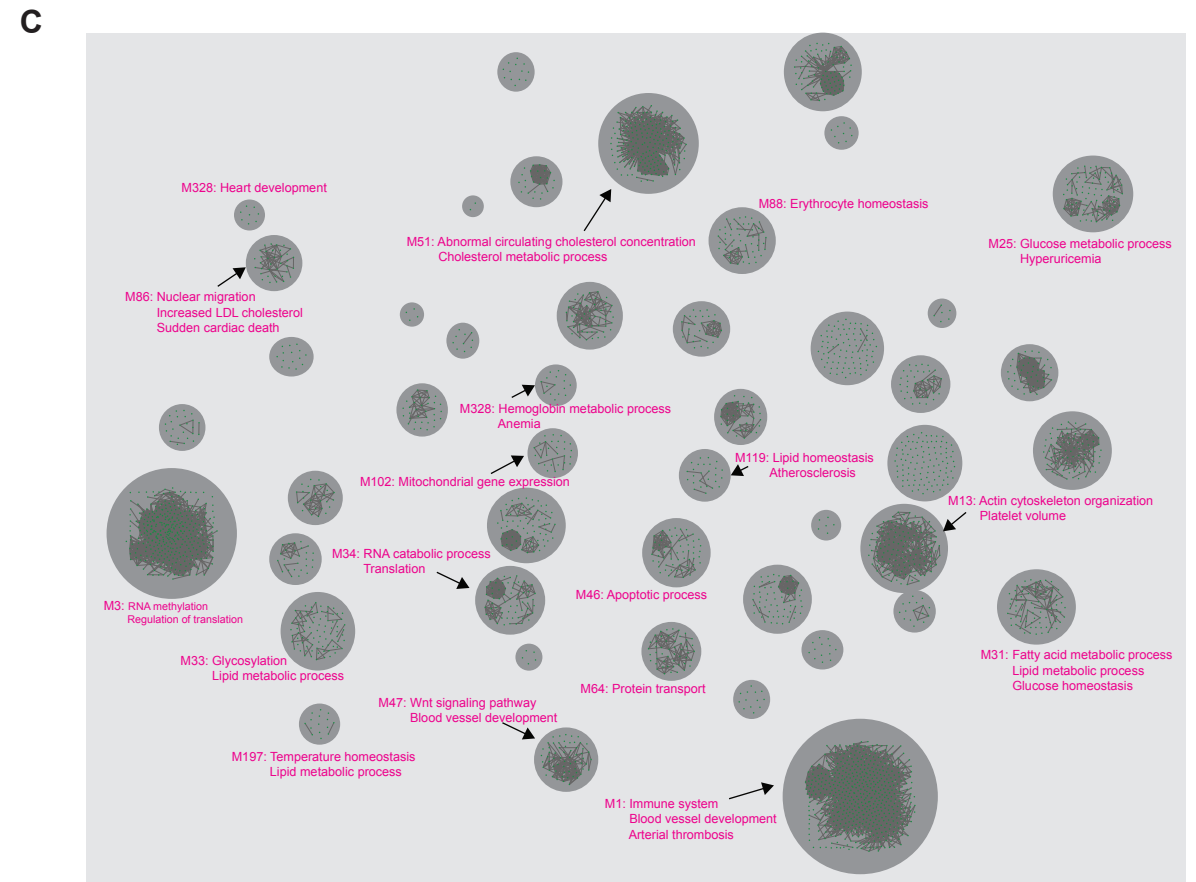
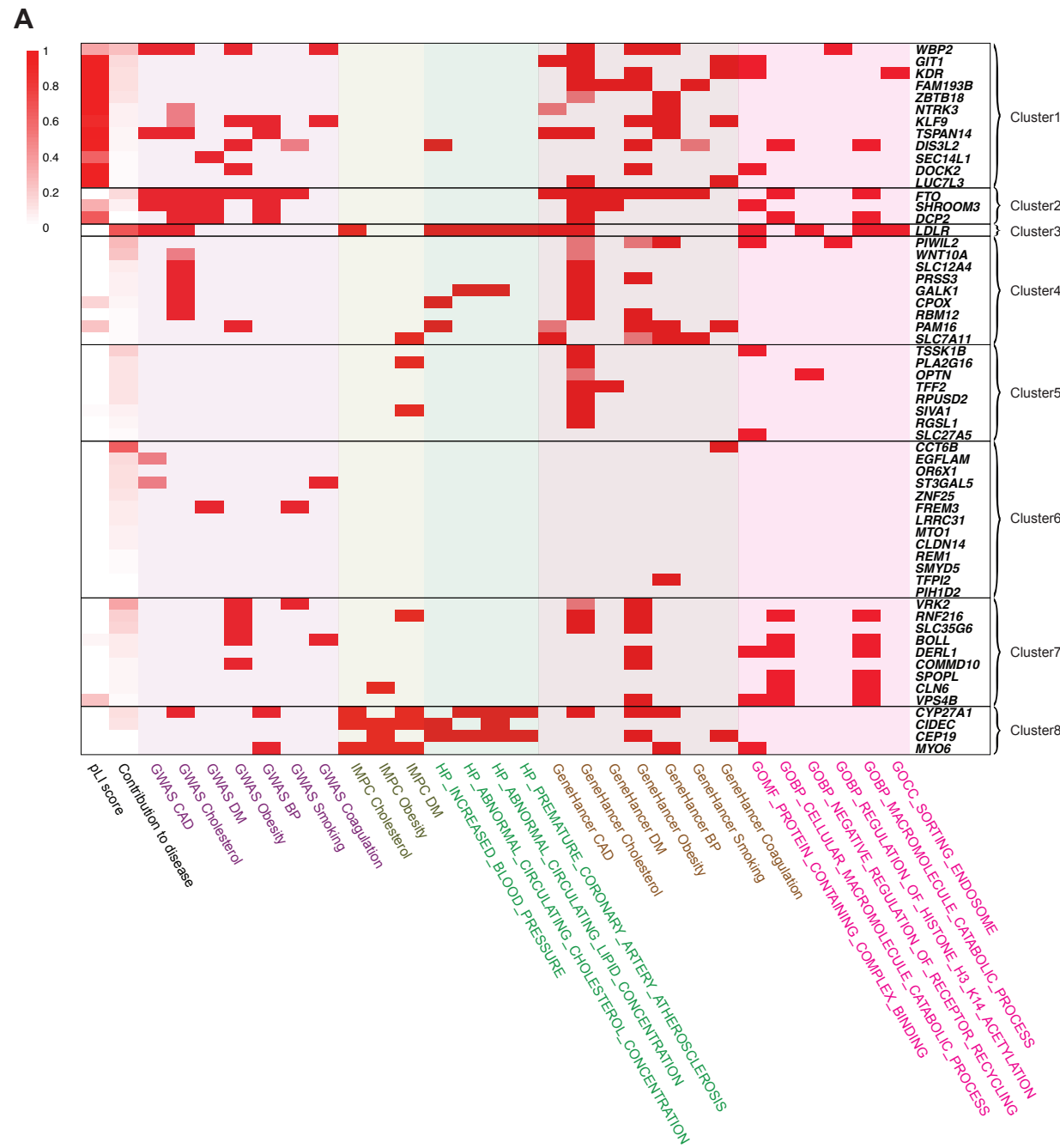
- 878 Vedantam, S., Buchkovich, M.L., Yang, J., et al. (2015). Genetic studies of body  
879 mass index yield new insights for obesity biology. *Nature* 518, 197–206.
- 880 33. Richardson, T.G., Sanderson, E., Palmer, T.M., Ala-Korpela, M., Ference, B.A.,  
881 Davey Smith, G., and Holmes, M.V. (2020). Evaluating the relationship between  
882 circulating lipoprotein lipids and apolipoproteins with risk of coronary heart  
883 disease: A multivariable Mendelian randomisation analysis. *PLoS Med.* 17,  
884 e1003062.
- 885 34. Sakaue, S., Kanai, M., Tanigawa, Y., Karjalainen, J., Kurki, M., Koshihara, S., Narita,  
886 A., Konuma, T., Yamamoto, K., Akiyama, M., et al. (2021). A cross-population  
887 atlas of genetic associations for 220 human phenotypes. *Nat. Genet.* 53, 1415–  
888 1424.
- 889 35. Zhuang, Z., Yao, M., Wong, J.Y.Y., Liu, Z., and Huang, T. (2021). Shared genetic  
890 etiology and causality between body fat percentage and cardiovascular diseases: a  
891 large-scale genome-wide cross-trait analysis. *BMC Med.* 19, 100.
- 892 36. Evangelou, E., Warren, H.R., Mosen-Ansorena, D., Mifsud, B., Pazoki, R., Gao, H.,  
893 Ntritsos, G., Dimou, N., Cabrera, C.P., Karaman, I., et al. (2018). Genetic analysis  
894 of over 1 million people identifies 535 new loci associated with blood pressure  
895 traits. *Nat. Genet.* 50, 1412–1425.
- 896 37. Sinnott-Armstrong, N., Tanigawa, Y., Amar, D., Mars, N., Benner, C., Aguirre, M.,  
897 Venkataraman, G.R., Wainberg, M., Ollila, H.M., Kiiskinen, T., et al. (2021).  
898 Genetics of 35 blood and urine biomarkers in the UK Biobank. *Nat. Genet.* 53,  
899 185–194.
- 900 38. Gargano, M.A., Matentzoglou, N., Coleman, B., Addo-Lartey, E.B.,  
901 Anagnostopoulos, A.V., Anderton, J., Avillach, P., Bagley, A.M., Bakštein, E.,  
902 Balhoff, J.P., et al. (2024). The Human Phenotype Ontology in 2024: phenotypes  
903 around the world. *Nucleic Acids Res.* 52, D1333–D1346.
- 904 39. Nordestgaard, B.G., Nicholls, S.J., Langsted, A., Ray, K.K., and Tybjaerg-Hansen,  
905 A. (2018). Advances in lipid-lowering therapy through gene-silencing technologies.  
906 *Nat. Rev. Cardiol.* 15, 261–272.
- 907 40. Raal, F.J., Rosenson, R.S., Reeskamp, L.F., Hovingh, G.K., Kastelein, J.J.P., Rubba,  
908 P., Ali, S., Banerjee, P., Chan, K.-C., Gipe, D.A., et al. (2020). Evinacumab for  
909 Homozygous Familial Hypercholesterolemia. *N. Engl. J. Med.* 383, 711–720.
- 910 41. Kessler, T., and Schunkert, H. (2021). Coronary Artery Disease Genetics  
911 Enlightened by Genome-Wide Association Studies. *JACC Basic Transl Sci* 6, 610–  
912 623.
- 913 42. Mortensen, M.B., and Nordestgaard, B.G. (2020). Elevated LDL cholesterol and  
914 increased risk of myocardial infarction and atherosclerotic cardiovascular disease  
915 in individuals aged 70–100 years: a contemporary primary prevention cohort.

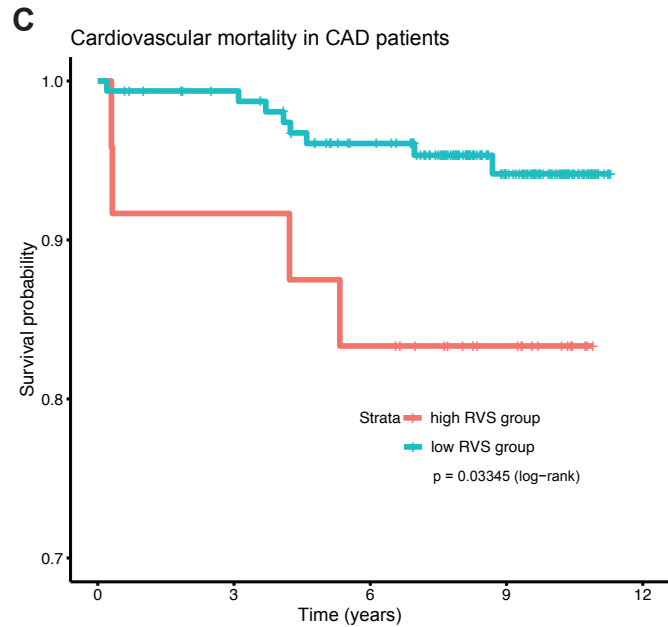
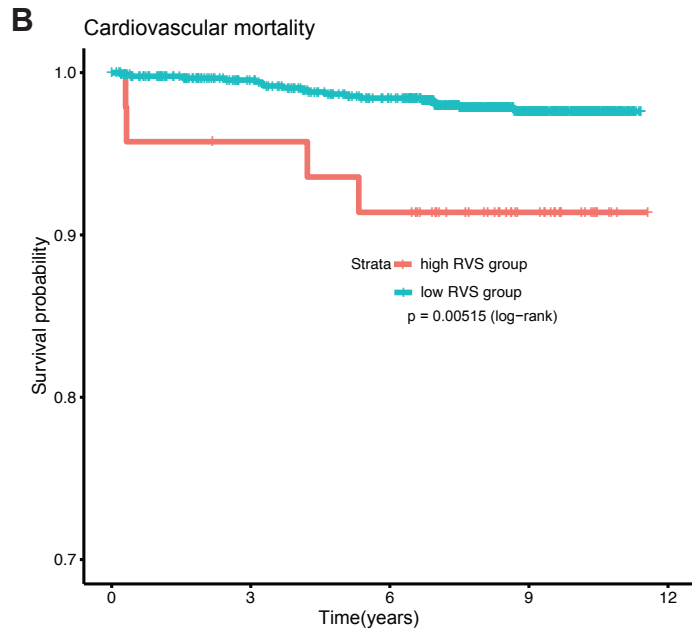
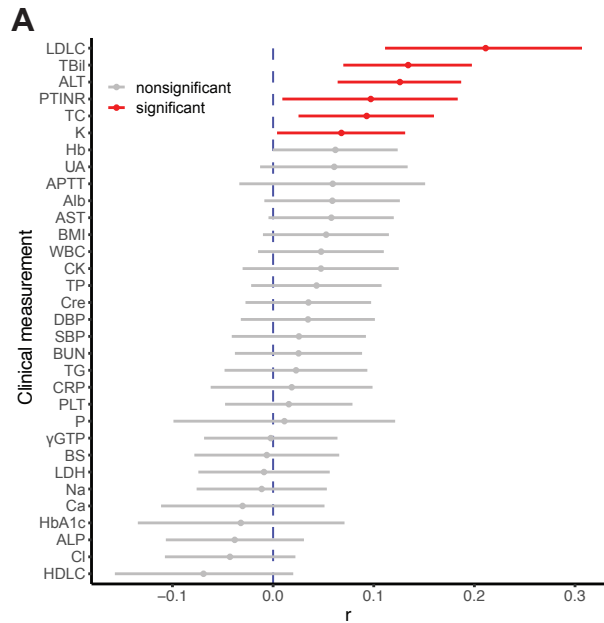
- 916           Lancet 396, 1644–1652.
- 917   43. Howard, B.V., Robbins, D.C., Sievers, M.L., Lee, E.T., Rhoades, D., Devereux,  
918       R.B., Cowan, L.D., Gray, R.S., Welty, T.K., Go, O.T., et al. (2000). LDL  
919       cholesterol as a strong predictor of coronary heart disease in diabetic individuals  
920       with insulin resistance and low LDL: The Strong Heart Study. *Arterioscler.*  
921       *Thromb. Vasc. Biol.* 20, 830–835.
- 922   44. Zhao, J.V., and Schooling, C.M. (2018). Coagulation Factors and the Risk of  
923       Ischemic Heart Disease. *Circulation: Genomic and Precision Medicine* 11,  
924       e001956.
- 925   45. Ndrepepa, G., and Kastrati, A. (2019). Alanine aminotransferase—a marker of  
926       cardiovascular risk at high and low activity levels. *J. Lab. Precis. Med.* 4, 29–29.
- 927   46. Shen, H., Zeng, C., Wu, X., Liu, S., and Chen, X. (2019). Prognostic value of total  
928       bilirubin in patients with acute myocardial infarction: A meta-analysis. *Medicine*  
929       98, e13920.
- 930   47. Emerging Risk Factors Collaboration, Di Angelantonio, E., Sarwar, N., Perry, P.,  
931       Kaptoge, S., Ray, K.K., Thompson, A., Wood, A.M., Lewington, S., Sattar, N., et al.  
932       (2009). Major lipids, apolipoproteins, and risk of vascular disease. *JAMA* 302,  
933       1993–2000.
- 934   48. Auer, P.L., and Lettre, G. (2015). Rare variant association studies: considerations,  
935       challenges and opportunities. *Genome Med.* 7, 16.
- 936   49. Chen, W., Coombes, B.J., and Larson, N.B. (2022). Recent advances and  
937       challenges of rare variant association analysis in the biobank sequencing era. *Front.*  
938       *Genet.* 13, 1014947.
- 939   50. Khetarpal, S.A., Babb, P.L., Zhao, W., Hancock-Cerutti, W.F., Brown, C.D., Rader,  
940       D.J., and Voight, B.F. (2018). Multiplexed Targeted Resequencing Identifies  
941       Coding and Regulatory Variation Underlying Phenotypic Extremes of High-  
942       Density Lipoprotein Cholesterol in Humans. *Circ Genom Precis Med* 11, e002070.
- 943   51. Diaz-Uriarte, R., Gómez de Lope, E., Giugno, R., Fröhlich, H., Nazarov, P.V.,  
944       Nepomuceno-Chamorro, I.A., Rauschenberger, A., and Glaab, E. (2022). Ten quick  
945       tips for biomarker discovery and validation analyses using machine learning. *PLoS*  
946       *Comput. Biol.* 18, e1010357.
- 947   52. Fahed, A.C., Wang, M., Homburger, J.R., Patel, A.P., Bick, A.G., Neben, C.L., Lai,  
948       C., Brockman, D., Philippakis, A., Ellinor, P.T., et al. (2020). Polygenic  
949       background modifies penetrance of monogenic variants for tier 1 genomic  
950       conditions. *Nat. Commun.* 11, 3635.
- 951   53. Chen, Z., and Schunkert, H. (2021). Genetics of coronary artery disease in the  
952       post-GWAS era. *J. Intern. Med.* 290, 980–992.

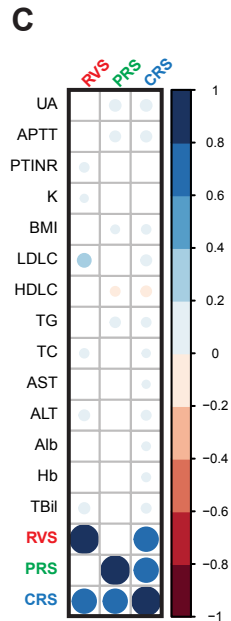
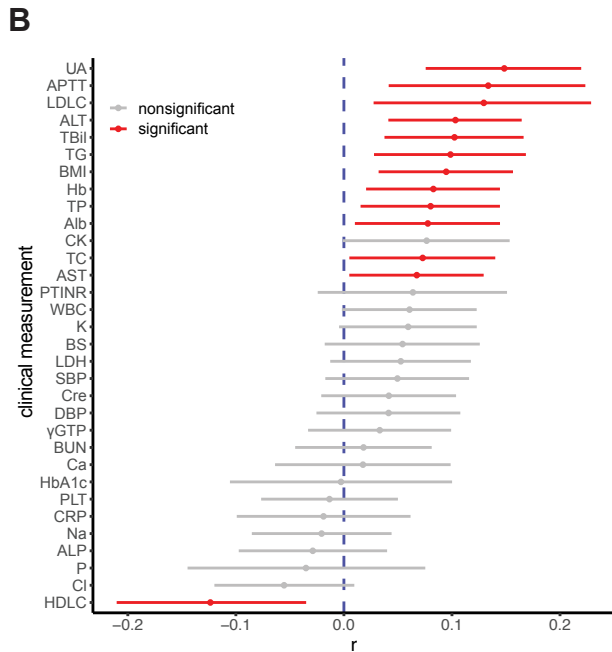
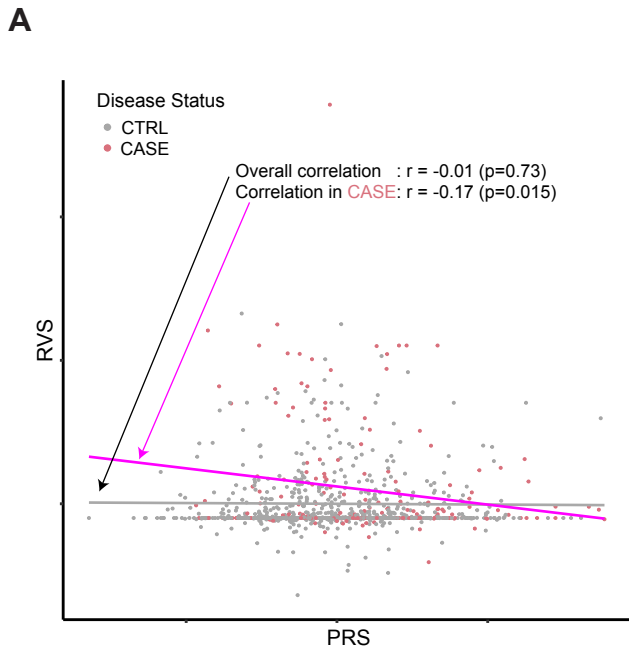
- 953 54. Martin, A.R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B.M., and Daly, M.J.  
954 (2019). Clinical use of current polygenic risk scores may exacerbate health  
955 disparities. *Nat. Genet.* 51, 584–591.
- 956 55. Nagai, A., Hirata, M., Kamatani, Y., Muto, K., Matsuda, K., Kiyohara, Y.,  
957 Ninomiya, T., Tamakoshi, A., Yamagata, Z., Mushirola, T., et al. (2017). Overview  
958 of the BioBank Japan Project: Study design and profile. *J. Epidemiol.* 27, S2–S8.
- 959 56. Hirata, M., Kamatani, Y., Nagai, A., Kiyohara, Y., Ninomiya, T., Tamakoshi, A.,  
960 Yamagata, Z., Kubo, M., Muto, K., Mushirola, T., et al. (2017). Cross-sectional  
961 analysis of BioBank Japan clinical data: A large cohort of 200,000 patients with 47  
962 common diseases. *J. Epidemiol.* 27, S9–S21.
- 963 57. Setoh, K., and Matsuda, F. (2022). Cohort Profile: The Nagahama Prospective  
964 Genome Cohort for Comprehensive Human Bioscience (The Nagahama Study). In  
965 *Socio-Life Science and the COVID-19 Outbreak: Public Health and Public Policy*,  
966 M. Yano, F. Matsuda, A. Sakuntabhai, and S. Hirota, eds. (Springer Singapore), pp.  
967 127–143.
- 968 58. Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with  
969 Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760.
- 970 59. Galinsky, K.J., Bhatia, G., Loh, P.-R., Georgiev, S., Mukherjee, S., Patterson, N.J.,  
971 and Price, A.L. (2016). Fast Principal-Component Analysis Reveals Convergent  
972 Evolution of ADH1B in Europe and East Asia. *Am. J. Hum. Genet.* 98, 456–472.
- 973 60. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J.  
974 (2015). Second-generation PLINK: rising to the challenge of larger and richer  
975 datasets. *Gigascience* 4, 7.
- 976 61. Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A., and  
977 Reich, D. (2006). Principal components analysis corrects for stratification in  
978 genome-wide association studies. *Nat. Genet.* 38, 904–909.
- 979 62. Zhou, W., Nielsen, J.B., Fritsche, L.G., Dey, R., Gabrielsen, M.E., Wolford, B.N.,  
980 LeFaive, J., VandeHaar, P., Gagliano, S.A., Gifford, A., et al. (2018). Efficiently  
981 controlling for case-control imbalance and sample relatedness in large-scale  
982 genetic association studies. *Nat. Genet.* 50, 1335–1341.
- 983 63. Zhou, W., Zhao, Z., Nielsen, J.B., Fritsche, L.G., LeFaive, J., Gagliano Taliun, S.A.,  
984 Bi, W., Gabrielsen, M.E., Daly, M.J., Neale, B.M., et al. (2020). Scalable  
985 generalized linear mixed model for region-based association tests in large biobanks  
986 and cohorts. *Nat. Genet.* 52, 634–639.
- 987 64. Zhou, W., Bi, W., Zhao, Z., Dey, K.K., Jagadeesh, K.A., Karczewski, K.J., Daly,  
988 M.J., Neale, B.M., and Lee, S. (2021). Set-based rare variant association tests for  
989 biobank scale sequencing data sets. *medRxiv*, 2021.07.12.21260400.
- 990 65. Ioannidis, N.M., Rothstein, J.H., Pejaver, V., Middha, S., McDonnell, S.K., Baheti,

- 991 S., Musolf, A., Li, Q., Holzinger, E., Karyadi, D., et al. (2016). REVEL: An  
992 Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants. *Am.*  
993 *J. Hum. Genet.* *99*, 877–885.
- 994 66. Wang, K., Li, M., and Hakonarson, H. (2010). ANNOVAR: functional annotation  
995 of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* *38*,  
996 e164.
- 997 67. A global reference for human genetic variation | Nature <https://www.nature.com> >  
998 articles<https://www.nature.com> > articles.
- 999 68. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q.,  
1000 Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The  
1001 mutational constraint spectrum quantified from variation in 141,456 humans.  
1002 *Nature* *581*, 434–443.
- 1003 69. Kolberg, L., Raudvere, U., Kuzmin, I., Adler, P., Vilo, J., and Peterson, H. (2023).  
1004 g:Profiler-interoperable web service for functional enrichment analysis and gene  
1005 identifier mapping (2023 update). *Nucleic Acids Res.* *51*, W207–W212.
- 1006 70. Blondel, V.D., Guillaume, J.-L., Lambiotte, R., and Lefebvre, E. (2008). Fast  
1007 unfolding of communities in large networks. *J. Stat. Mech.* *2008*, P10008.
- 1008 71. Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin,  
1009 N., Schwikowski, B., and Ideker, T. (2003). Cytoscape: a software environment for  
1010 integrated models of biomolecular interaction networks. *Genome Res.* *13*, 2498–  
1011 2504.
- 1012 72. Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I.,  
1013 Chew, E.Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype  
1014 imputation service and methods. *Nat. Genet.* *48*, 1284–1287.

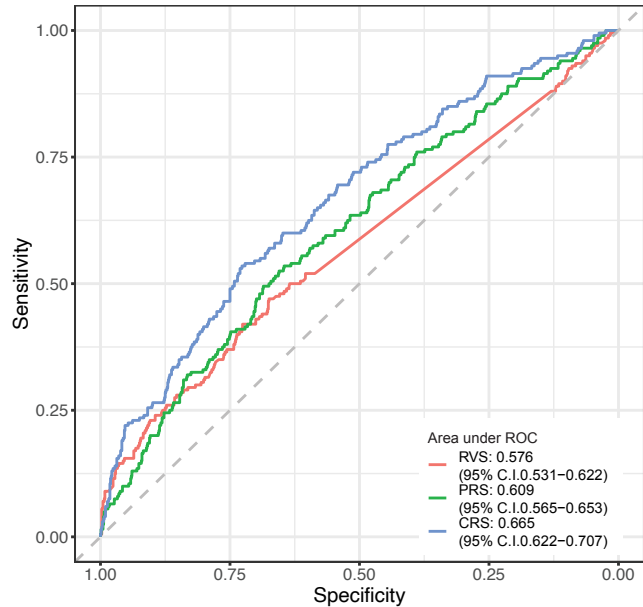




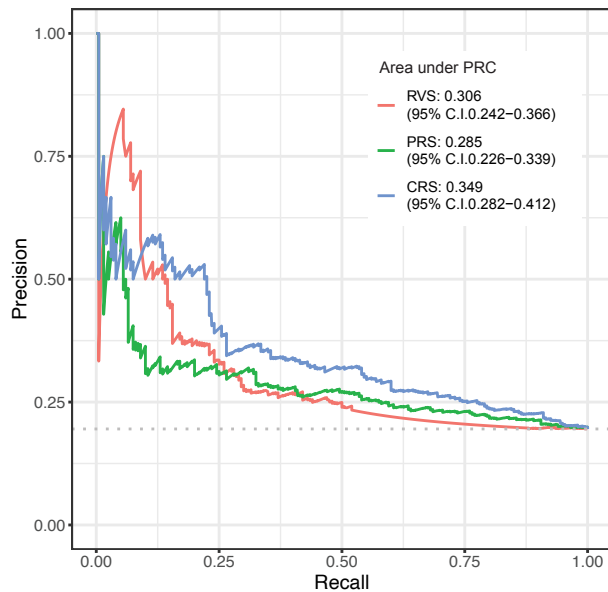






**A** ROC analysis**B**

P-R curve

**C**