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 ^{1,2,3,4} , Kaoru Ito ^{1*} , Sai Zhang ⁵ , Satoshi Koyama ^{1,6,7} , Martin Kjellberg ^{3,8} ,
6	Hiroki Yoshida ^{1,2} , Ryo Kurosawa ¹ , Hiroshi Matsunaga ^{1,2} , Kazuo Miyazawa ¹ ,
7	Nobuyuki Enzan ^{1,6,7} , Changhoon Kim ⁹ , Jeong-Sun Seo ^{9,10} , Koichiro Higasa ^{11,12} ,
8	Kouichi Ozaki ^{1,13} , Yoshihiro Onouchi ^{1,14} , The Biobank Japan Project, Koichi Matsuda ¹⁵ ,
9	Yoichiro Kamatani ¹⁶ , Chikashi Terao ¹⁷ , Fumihiko Matsuda ¹² , Michael Snyder ^{3,4*} ,
10	Issei Komuro ^{18,19*}
11	
12	Affiliations:
13	¹ Laboratory for Cardiovascular Genomics and Informatics, RIKEN Center for
14	Integrative Medical Sciences, Yokohama, Japan.
15	² Department of Cardiovascular Medicine, Graduate School of Medicine, The
16	University of Tokyo, Tokyo, Japan.
17	³ Department of Genetics, Center for Genomics and Personalized Medicine, Stanford
18	University School of Medicine, Stanford, USA.
19	⁴ Stanford Cardiovascular Institute, Stanford University School of Medicine, Stanford,
20	USA.
21	⁵ Department of Epidemiology, University of Florida, Gainesville, USA.
22	⁶ Center for Genomic Medicine, Department of Medicine, Massachusetts General
23	Hospital, Boston, MA, USA.
24	⁷ Program in Medical and Population Genetics, Broad Institute of Harvard and MIT,
	1

- 25 Cambridge, MA, USA.
- ⁸ School of Engineering Sciences in Chemistry, Biotechnology and Health, Royal
- 27 Institute of Technology (KTH), Stockholm, Sweden
- ⁹ Bioinformatics Institute, Macrogen Inc., Seoul, Republic of Korea.
- ¹⁰ Asian Genome Institute, Seoul National University Bundang Hospital, Gyeonggi-do,
- 30 Republic of Korea.
- 31 ¹¹ Department of Genome Analysis, Institute of Biomedical Science, Kansai Medical
- 32 University, Hirakata, Japan.
- 33 ¹² Human Disease Genomics, Center for Genomic Medicine, Kyoto University
- 34 Graduate School of Medicine, Kyoto, Japan.
- 35 ¹³ Medical Genome Center, Research Institute, National Center for Geriatrics and
- 36 Gerontology, Obu, Japan.
- ¹⁴ Department of Public Health, Chiba University Graduate School of Medicine, Chiba,

38 Japan.

- ¹⁵ Department of Computational Biology and Medical Science, Graduate School of
- 40 Frontier Sciences, The University of Tokyo, Tokyo, Japan.
- 41 ¹⁶Laboratory of Complex Trait Genomics, Department of Computational Biology and
- 42 Medical Sciences, Graduate School of Frontier Sciences, The University of Tokyo,
- 43 Tokyo, Japan.
- ¹⁷ Laboratory for Statistical and Translational Genetics, RIKEN Center for Integrative
- 45 Medical Science, Yokohama, Japan
- 46 ¹⁸ International University of Health and Welfare, Tokyo, Japan
- 47 ¹⁹ Department of Frontier Cardiovascular Science, Graduate School of Medicine, The
- 48 University of Tokyo, Tokyo, Japan

49

50

- 52 kaoru.ito@riken.jp (K.I.)
- 53 mpsnyder@stanford.edu (M.S.)
- 54 komuro-tky@umin.ac.jp (I.K.)

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

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

76

77

79 Introduction

80 Despite advancements in treatments and medications, coronary artery disease 81 (CAD), encompassing conditions such as angina pectoris and myocardial infarction (MI), remains a leading cause of death worldwide ^{1,2}. CAD etiology is complex, 82 83 involving a multifaceted interplay between genetic predisposition and environmental 84 determinants. Lifestyle factors including diet, smoking, and physical activity are wellestablished contributors to the onset and progression of CAD^{3,4}. Additionally, 85 conditions such as elevated low-density lipoprotein (LDL) cholesterol, hypertension, 86 and glucose intolerance further exacerbate the risk profile ⁵. The importance of genetic 87 88 predisposition is also underscored by a European twin study, which estimated that genetic factors contributed to over 50% of CAD development ^{6,7}. Therefore, 89 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 CAD⁸⁻¹². Polygenic risk scores 95 (PRSs) derived from GWAS summary statistics have enabled the estimation of individual-level CAD risk ^{13,14}. However, despite these significant advancements, the 96 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 variants, while rare variants are often underrepresented in these analyses ^{5,15}. Rare 99 100 variants often have a large effect size on diseases and phenotypes, making them a promising target for drug development ¹⁶. Incorporating rare variants into genetic risk 101 102 scores could significantly enhance the accuracy of CAD prediction. Despite this

potential, previous GWASs and aggregated rare variant association analyses have struggled even in large-scale sequencing studies, identifying only a few genes at exomewide significance per trait ^{17,18}. Furthermore, calculating a genetic risk score based on rare variants is challenging because gene-level effect sizes are not estimated by 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 aortic aneurysm¹⁹. In the current study, we conducted whole genome sequencing 113 114 (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 3,148 controls. In order to enhance the genetic discovery power ²⁰, we prioritized 127 patients with early-onset MI, a severe form of CAD, from the BioBank Japan (BBJ) 128 129 cohort. The average age of MI onset in these patients was 47.4 ± 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**).

We first conducted a single variant association test in the discovery cohort using a logistic regression model implemented in PLINK software with covariates of age, sex and top ten ancestry principal components (PCs). The genomic inflation factor λ_{GC} was calculated to be 1.03, indicating minimal inflation of test statistics and suggesting that the quality control applied to the samples was adequate (**Figure S2**). This initial single variant association analysis did not identify any genetic loci that 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**))^{10,11}. 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 (SKAT-O). While no genomic inflation was observed ($\lambda = 0.939$) (Figure S4), the 159 LDLR gene surpassed a suggestive threshold ($p = 2.3 \times 10^{-5}$). However, no genes reached 160 the gene-wide significance threshold of $p = 2.5 \times 10^{-6}$ (Figure S4 and Table S4). This 161 162 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

The machine learning-based framework prioritizes disease-associated genes and reveals molecular networks

We next conducted a machine learning-based rare variant analysis using a modified HEAL ¹⁹. In this framework, we first quantified the mutation burden for each gene in each participant defined by the cumulative effects of deleterious 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_{CAD} genes, we assessed 177 constraint scores and checked for overlaps with neighboring genes identified in previous GWASs on CAD and its risk factors. Using the Genehancer database²¹, which provides 178 179 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) ²² database to investigate the phenotypes associated with a gene knockout (KO) in mice 182 183 and conducted gene set enrichment analysis to identify functional clusters among the 184 HEAL_{CAD} 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 familial hypercholesterolemia²³ and has been consistently associated with CAD in 189 previous GWASs and genome sequencing studies ^{9,24,25}, supporting the validity of our 190 191 machine learning-based framework. In the IMPC database, LDLR KO mice showed increased circulating cholesterol levels ²⁶, a known risk factor for CAD. Cluster 7 192 193 contained genes related to obesity and metabolic processes, such as the RNF216 locus, which is associated with body mass index (BMI)²⁷ and increased glucose levels in KO 194

mice 22 . Additionally, the VRK2 locus has been reported to be associated with BMI 28 , 195 smoking behavior and alcohol use ²⁹, indicating its broader impact on metabolic health. 196 197 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 highlighted for its strong association with obesity ^{30,31} and related phenotypes linked to 199 BMI ³², LDL cholesterol ³³, blood pressure ³⁴, and CAD ³⁵. Cluster 8 encompassed 200 201 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, the CYP27A1 locus is associated with diastolic blood pressure ³⁶ and triglyceride levels 203 ³⁷ and has connections to cholesterol levels and premature CAD according to human 204 phenotype ontology 38 . 205

To further determine the functions of the fifty-nine genes, we mapped them 206 207 onto the human protein-protein interaction (PPI) network followed by identifying proteins that were tightly clustered with these HEAL_{CAD} genes as topological modules 208 209 ¹⁹. We identified 46 tightly clustered topological modules encompassing the HEAL_{CAD} 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 2.53×10^{-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 ANGPTL3 (Figure 2D), which are well-known targets for medications treating 217 dyslipidemia and CAD ^{39 40}. Other modules were associated with different biological 218

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 41 , 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.

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

potassium levels (**Figure 3A and Table S9**). These correlations are noteworthy since elevated cholesterol levels and coagulation abnormalities are established risk factors for CAD $^{42-44}$. Moreover, alterations of total bilirubin and AST were also reported to be associated with cardiovascular risk 45,46 , reinforcing the clinical relevance of the RVS in 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

The integration of RVS and PRS improves the performance of the genomic risk score

259 Many GWASs have been conducted for CAD, leading to the development of PRS that primarily comprise common variants to predict the risk of CAD. Multiple 260 261 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 CAD-GWAS¹¹ in the validation cohort. The PRS also significantly predicted CAD with 266

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 = 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 against CAD (Figure 3A, Figure S7 and Table S10)⁴⁷. These data support the notion 276 277 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 urinary 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 validation cohort (AUROC 0.66 vs 0.61, p=0.007; Pseudo R^2 0.093 vs 0.040, p = 289 0.0018; AUPRC 0.35 vs 0.29, p = 0.0154) (Figure 5 and Table S12). These results 290

- 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
- addressing some of the unexplained heritability in the disease.

294

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 ^{48,49}. While some studies have adopted a targeted resequencing approach by selecting specific genes based on prior knowledge ^{25,50}; previous attempts 312 313 at genome-wide or exome-wide analyses have often suffered from insufficient statistical 314 power, leading to limited success in identifying previously uncharacterized genes associated with complex traits like CAD ²⁰. Also in this study, the single variant 315 316 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 single gene, LDLR, reached a significance level in the gene-based test for CAD¹⁷. 319

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 aneurysm¹⁹. We adapted and optimized the model for CAD patients, marking the first 324 325 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 crossvalidation and an independent validation cohort ⁵¹. 333

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 genetic mutations with PRS have been reported ⁵², our study presented a more 357 358 comprehensive approach to combine rare and common variant information. Furthermore, these findings reinforce the recognition that rare variants, despite their low 359 360 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.

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

368 study have unknown functions, especially in cluster 6. However, many loci and genes identified in GWAS on CAD remain functionally uncharacterized, as well ^{41,53}. 369 370 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 less accurate in other populations ^{11,54}. These results need to be validated in other 374 375 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.

385 Consortia

386 The Biobank Japan Project

- 387 Koichi Matsuda^{1,2}, Takayuki Morisaki^{2,3}, Yukinori Okada⁴, Yoichiro Kamatani⁵, Kaori
- 388 Muto⁶, Akiko Nagai⁶, Yoji Sagiya², Natsuhiko Kumasaka⁷, Yoichi Furukawa⁸, Yuji
- 389 Yamanashi³, Yoshinori Murakami³, Yusuke Nakamura³, Wataru Obara⁹, Ken Yamaji¹⁰,
- 390 Kazuhisa Takahashi¹¹, Satoshi Asai^{12,13}, Yasuo Takahashi¹³, Shinichi Higashiue¹⁴, Shuzo
- 391 Kobayashi¹⁴, Hiroki Yamaguchi¹⁵, Yasunobu Nagata¹⁵, Satoshi Wakita¹⁵, Chikako Nito¹⁶,
- 392 Yu-ki Iwasaki¹⁷, Shigeo Murayama¹⁸, Kozo Yoshimori¹⁹, Yoshio Miki²⁰, Daisuke
- 393 Obata²¹, Masahiko Higashiyama²², Akihide Masumoto²³, Yoshinobu Koga²³ & Yukihiro
- 394 Koretsune²⁴
- 395
- ^{1.}Laboratory of Genome Technology, Human Genome Center, Institute of Medical
 Science, The University of Tokyo, Tokyo, Japan.
- ² Laboratory of Clinical Genome Sequencing, Graduate School of Frontier Sciences,
- 399 The University of Tokyo, Tokyo, Japan.
- 400 ³ The Institute of Medical Science, The University of Tokyo, Tokyo, Japan.
- 401 ⁴ Department of Genome Informatics, Graduate School of Medicine, The University of
- 402 Tokyo, Tokyo, Japan.
- 403 ⁵ Laboratory of Complex Trait Genomics, Graduate School of Frontier Sciences, The
- 404 University of Tokyo, Tokyo, Japan.
- ⁶ Department of Public Policy, Institute of Medical Science, The University of Tokyo,
 Tokyo, Japan.
- ⁷ Division of Digital Genomics, Institute of Medical Science, The University of Tokyo,
- 408 Tokyo, Japan.

- 409 ⁸ Division of Clinical Genome Research, Institute of Medical Science, The University of
- 410 Tokyo, Tokyo, Japan.
- ⁹ Department of Urology, Iwate Medical University, Iwate, Japan.
- 412 ¹⁰ Department of Internal Medicine and Rheumatology, Juntendo University Graduate
- 413 School of Medicine, Tokyo, Japan.
- 414 ¹¹ Department of Respiratory Medicine, Juntendo University Graduate School of
- 415 Medicine, Tokyo, Japan.
- 416 ¹² Division of Pharmacology, Department of Biomedical Science, Nihon University
- 417 School of Medicine, Tokyo, Japan.
- 418 ¹³ Division of Genomic Epidemiology and Clinical Trials, Clinical Trials Research
- 419 Center, Nihon University. School of Medicine, Tokyo, Japan.
- 420 ¹⁴ Tokushukai Group, Tokyo, Japan.
- 421 ¹⁵ Department of Hematology, Nippon Medical School, Tokyo, Japan.
- 422 ¹⁶ Laboratory for Clinical Research, Collaborative Research Center, Nippon Medical
- 423 School, Tokyo, Japan.
- 424 ¹⁷ Department of Cardiovascular Medicine, Nippon Medical School, Tokyo, Japan.
- ¹⁸ Tokyo Metropolitan Geriatric Hospital and Institute of Gerontology, Tokyo, Japan.
- 426 ¹⁹ Fukujuji Hospital, Japan Anti-Tuberculosis Association, Tokyo, Japan.
- 427 ²⁰ The Cancer Institute Hospital of the Japanese Foundation for Cancer Research, Tokyo,
- 428 Japan.
- 429 ²¹ Center for Clinical Research and Advanced Medicine, Shiga University of Medical
- 430 Science, Shiga, Japan.
- 431 ²² Department of General Thoracic Surgery, Osaka International Cancer Institute, Osaka,
- 432 Japan.

433 ²³ Iizuka Hospital, Fukuoka, Japan.

434 ²⁴ 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

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

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 _{CAD} 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

500 Figure 2. Functional analysis of HEAL_{CAD} 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 $HEAL_{CAD}$ 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.

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 Pearson's correlation coefficients and their 95% confidence intervals (CIs). Exact P 527 528 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 protein; Cre, creatinine; DBP, diastolic blood pressure; SBP, systolic blood pressure; 538 539 BUN, blood urea nitrogen; TG, triglycerides; CRP, C-reactive protein; PLT, platelet; P, 540 Phosphorus; yGTP, gamma-glutamyl transpeptidase; BS, blood sugar; LDH, Lactate 541 dehydrogenase.

542

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	γ 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 Risk Score). The curve plots the true positive rate (sensitivity) against the false positive 565 566 rate (1-specificity) for different threshold values of the predictive score. The area under the curve (AUC) is indicated, representing the score's accuracy in predicting the 567 568 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 20,000 times bootstrap replication method. (C) Boxplot of Pseudo R² for CAD 573 prediction performance. This box plot displays the pseudo- R^2 values comparing the 574 CAD prediction performance of RVS, PRS and CRS. The distribution of pseudo- R^2 was 575 576 estimated from 20,000 times bootstrapping. The box plot center line represents the median, the bounds represent the first and third quartile, and the whiskers reach to 1.5 577 578 times the interquartile range.

579 Tables

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

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

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

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

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

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 genetic contribution to disease²⁰, we prioritized patients with early-onset MI, one of the 600 601 most severe forms of CAD, within the BBJ cohort for WGS (age of MI onset in 15x and 602 30x WGS cohort: 47.4 ± 4.1 years and 36.0 ± 3.9 years, respectively). Sequenced reads were aligned to the hs37d5 reference genome using BWA software ⁵⁸. The genotypes of 603 604 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.

(1) For 15x depth data, filtered depth (DP) < 2, quality of the assigned genotype 607 608 (genotype quality; GQ) < 20. (2) For 30x depth data, DP < 5, GQ < 20, DP > 60 and GQ609 < 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 (Pi-hat > 0.1875) and (5) excess heterozygosity. We also excluded non-Japanese 615 participants estimated from Principal component analysis (PCA) calculated using PLINK 2.0⁵⁹. The total number of genomes that failed data quality control is 616 617 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) Hardy-Weinberg equilibrium (P < 1 $*10^{-6}$), (3) variants in the low complexity region. 621 622 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

The single variant association test was performed by logistic regression implemented in PLINK 2.0⁶⁰ with adjustment for age, sex, and the first 10 principal components of ancestry. Principal components of ancestry were calculated using PLINK 2.0⁵⁹. The inclusion of principal components as covariates in the logistic regression analysis increases the power to detect true genetic associations and minimizes

confounding by population stratification ⁶¹. Variants with a missing rate of less than 631 0.01 were included in the analysis. Genomic inflation factor (λ_{GC}) was calculated using 632 variants with MAF \geq 0.001. Single variant association analysis was also performed 633 using SAIGE ⁶² with adjustment for age, sex, and the first 10 principal components of 634 635 ancestry. SAIGE is widely used in GWASs for binary traits to account for population structure and relatedness while correcting for the type I error rates ⁶². The genome-wide 636 significance threshold was set at $P = 5 * 10^{-8}$. To define a locus, we added 500 kb to 637 638 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+ software, which accounts for the relatedness among the study samples ^{63,64}. We first 645 646 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, frameshift deletion, stopgain, stoploss, and splice site variants). Splice-site variants, 650 pLOF variants and damaging missense variants defined by a REVEL score $> 0.5^{65}$ were 651 652 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. Gene-wide significance threshold and suggestive threshold were set at $P = 2.5 * 10^{-6}$ 654

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 HEAL method is described in the original paper¹⁹. In this framework (Figure S8), we 660 first annotated each variant using ANNOVAR software ⁶⁶ and extracted rare 661 nonsynonymous variants (nonsynonymous SNV, nonframeshift insertion, nonframeshift 662 663 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 Genomes Project ⁶⁷. Variants with high frequency in the WGS data and gnomAD East 665 Asian database 68 (MAF > 0.1) were also filtered. To estimate the mutation burden for 666 667 each gene based on the rare variants, we used the REVEL score (ranges from 0 to 1 with a higher score indicating a damaging variant), which was internally computed by 668 ANNOVAR software. The deleteriousness score of the putative loss of function (pLOF) 669 670 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 ($x_n = (g_{1n}, g_{2n}, ..., g_{mn})$, where m is the number of the total genes). The 676 mutation burden was standardized (Z-score normalization). We trained a regularized

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

$$\widehat{y_n} = P(y_n = 1 | \boldsymbol{x_n}) = \sigma(\boldsymbol{w}^T \boldsymbol{x_n}) = \frac{1}{1 + exp(-\boldsymbol{w}^T \boldsymbol{x_n})}$$

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

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

685 In this regularized logistic regression, regularization strength is determined by 686 parameter λ , 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.

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 the constraint score (pLI) from the gnomAD database 68 , identifying whether the genes 701 702 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 GeneHancer database, which includes genome-wide enhancers and their target genes ²¹. 706 707 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 performed using g:Profiler ⁶⁹ to gain insights into the biological processes and human 711 phenotypic abnormalities associated with these genes ²². We considered statistical 712 713 significance for the enrichment analysis with a false discovery rate under 0.1.

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

the network into different modules. Following algorithm convergence, we obtained 1,261 modules with an average size of 13 nodes. Among the 1,261 PPI modules, 46 encompassed at least one gene identified by the machine learning analysis. We used g:profiler to determine functional enrichment for each module. Cytoscape software ⁷¹ 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) Hardy–Weinberg equilibrium (P < 1 $*10^{-6}$) and (3) heterozygous counts <5 were 742 743 excluded. We performed pre-phasing using Eagle software. Phased haplotypes were imputed to the in-house reference panel from BBJ ¹¹ by minimac3 ⁷². Variants with low 744

imputation quality ($R^2 < 0.3$) were excluded. GWAS was performed by logistic regression implemented in PLINK 2.0 ⁶⁰ with adjustment for age, age², sex and first 10 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, $a_{i,i}$ is the number of effect allele of *j*th 749 variant in *i*th sample, and β_i 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 thresholding method ¹³. The relationship between RVS and PRS was examined by 751 Pearson's correlation coefficient, both in cases only and across the validation cohort. We 752 753 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) and Nagelkerke's pseudo R^2 metrics. The P values were derived using a 20000 times 760 761 bootstrap replication method. In all statistical analyses, R software was used and a two-762 sided P < 0.05 was considered statistically significant.

763

765 **Reference**

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

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

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

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

916 Lancet *396*, 1644–1652.

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

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

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









Α







С

A ROC analysis





