1 # Title

2 ## **Development and validation of a next-gen health stratification**

3 ## **engine to determine risk for multiple cardiovascular diseases**

4

5 Mehrdad Rezaee [1,2*], Arsia Takeh[1], Igor Putrenko[1], Andrea Ganna[3,4,5], and Erik Ingelsson[6,7,8]

6

7 [1] Precision Wellness Inc., 1901 Embarcadero Rd #102, Palo Alto, CA, USA;

8 [2] Cardiac and Vascular Care, Inc. 2030 Forest Ave, San Jose, CA, USA;

9 [3] Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge,

10 MA, USA;

11 [4] Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA,

12 USA;

13 [5] Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General

14 Hospital, Boston, MA, USA;

15 [6] Department of Medicine, Division of Cardiovascular Medicine, Stanford University School of

16 Medicine, Stanford, CA 94305, USA;

17 [7] Stanford Cardiovascular Institute, Stanford, CA 94305, USA;

18 [8] Stanford Diabetes Research Center, Stanford, CA 94305, USA.

19

20 *Corresponding Author

21 E-mail: mrezaee@precisionwellness.com (MR)

22

23

2

## Abstract

Cardiometabolic diseases (CMD) impose greater impact on every aspect of health care than any other disease group. Accurate and in-time risk assessment of individuals for their propensity to develop CMD events is one of the most critical paths in preventing these conditions. The principal objective of the present study is to report the development, and validation of a next generation risk engine to predict CMD. UK Biobank population data was used to derive predictive models for six CMD. Missing data were imputed using imputation algorithms. Cox proportional hazard models were used to estimate annual absolute risk and relative risk of different risk factors for these conditions. In addition to conventional risk factors, the applied model included socioeconomic data, lifestyle factors and comorbidities as predictors of outcomes. In total, 416,936 individuals were included in the analysis. The derived prediction models achieved consistent and moderate-to-high discrimination performance (C-index) for all diseases: coronary artery disease (0.79), hypertension (0.82), type 2 diabetes mellitus (0.87), stroke (0.79), deep vein thrombosis (0.75), and abdominal aortic aneurysm (0.90). These results were consistent across age groups (37-73 years) and showed similar predictive abilities amongst those with pre-existing diabetes or hypertension. Calibration of risk scores showed that there was moderate overestimation of CMD-related conditions only in the highest decile of risk scores for all models. In summary, the newly developed algorithms, based on Cox proportional models, resulted in high disclination and good calibration for several CMD. The integrations of these algorithms on a single platform may have direct clinical impact.

# Introduction

Cardiometabolic diseases (CMD) continue to be the leading causes of death in the United States since the 1920s, and 45% of the U.S. population is projected to suffer from any of these diseases by 2035 [1]. The healthcare cost associated with these diseases represent one of the greatest global economic burdens [2]. As with any chronic condition, appropriate prevention and selective treatment for CMD are the most effective approaches to defer their clinical and financial impact on individuals and across populations.

Primary prevention of chronic diseases is a resource intensive, costly, and non-effective if applied through non-selective implementation [3]. Therefore, accurate population and individual stratification is needed to provide individualized, as well as population-specific care. In order to achieve clinically relevant risk stratification, established risk factors and novel population-specific data should be considered to derive clinically applicable prediction algorithms.

For over 20 years, the concept of cardiovascular risk assessment has been tested through prediction models that are utilized in the clinical setting [4-6]. Current prediction models have good discrimination abilities to identify individuals who will develop CMD. However, there are opportunities to address the limitations of current models, such as inclusion of contemporary risk factors, biomarkers and genetic information as part of the algorithms [7]. Also, the currently systems are limited to only a few diseases, such as coronary artery disease and stroke, without consideration of major comorbidities. Moreover, current models do not allow for imputation for missing data; and finally, they are primarily directed to prevention of disease over a 10-year span. In this study, the development and validation of a next-gen stratification platform that integrates conventional clinical risk factors and biomarkers,

69    socioeconomic, lifestyle factors and other co-morbidities data for six cardiometabolic diseases

70    (CMD) is presented.  To derive these new predictions models, we used data provided by the

71    UK Biobank (UKBB) project [8], including over 400,000 men and women aged 37–73 years,

72    with 6.1 years of median longitudinal follow-up.

73

5

# Materials and methods

## Baseline data preparation

Baseline data on 502,616 UKBB participants collected at assessment centers to derive the prediction models. Overall, 95% of the UKBB participants were self-described as white, with women comprising 54.4% of the total. CMD outcomes were determined based on International Classification of Diseases (ICD) edition 10 (ICD-10) codes, as well as self-reports for coronary artery disease (CAD), hypertension (HPT), type 2 diabetes mellitus (DM2), and deep vein thrombosis (DVT), and medications for CAD, HPT, and DM2. Six distinct datasets for each CMD were engineered. CAD was defined as I20–I25 and T82 codes. HPT was defined as I10, I15, and R03.0 codes. DM2 was defined as E11, E13, and E14 codes. Stroke was defined as G46.3, G46.4, I63, I66, I67, and I693 codes. DVT was defined as H34.8, H40.8, I23.6, I24.0, I63, I67.6, I74, I81, I82, I87.2, I87.3, K64.5, N48.8, N52.0, O03.3, O03.8, O04.8, O07.3, O08.7, 022, O87, Q26, T82.8, T83.8, T84.8, T85.8, and Z86.7 codes. Abdominal aortic aneurysm (AAA) was defined as I71 and I79.0 codes.

The UKBB data were subsequently linked to hospital episode statistics (HES) data from hospitals in England, Scotland and Wales. The age and date of a CMD event were determined based on primary or secondary ICD-10 codes in the HES data corresponding to the event using the earliest hospital record. The date of inclusion into the UKBB was defined as baseline and was used as starting point for time-to-event calculations. The exit date was determined as either date of death, end of follow-up (February 29, 2016), or a CMD event, whichever happened first. Only those CMD-positive cases that were identified by ICD-10 codes, self-reports, or medication as described above and had the date of the event determined

96    based on the HES data were included into analyses, reducing the number of participants to

97    416,936. In addition, participants with prior CMD events (before baseline) were excluded

98    from analyses of that specific event, e.g. those with prior CAD event were excluded from the

99    CAD analyses and so on.

100   The datasets created for each CMD were spitted into training and testing sets based on

101   80%/20% ratio. Testing sets were used for model validation and calibration. Age- and CMD-

102   specific testing sets were created by applying corresponding age and disease filters onto

103   general test datasets (without reusing any data from the training sets to avoid overfitting).

## Variable definition

104

105   To develop highly predictive CMD risk prediction models, in addition to using already

106   available UKBB data fields, the new variables were derived that captured sociodemographic

107   and socioeconomic factors, laboratory test results, physiological measurements, physical

108   activity, nutrition, alcohol consumption, family history of CMD; as well as the presence of

109   diseases, disorders, or previous surgeries as shown in Table 1.

110

111

112

113

114

115

116

117

118 **Table 1. Profile of variables for predicting the risk of six CMD.**

| | Type | N | 5% | 25% | 50% | 75% | 95% |
|---|---|---|---|---|---|---|---|
| **Sex** | binary | 416936 | women | women | women | men | men |
| **BMI** | continuous | 414265 | 20.84 | 23.94 | 26.48 | 29.58 | 35.82 |
| **DBP** | continuous | 415742 | 66 | 75 | 81 | 88 | 98.5 |
| **Age** | continuous | 416936 | 42 | 49 | 57 | 63 | 68 |
| **FEV1** | continuous | 376770 | 60.83 | 82.19 | 93.79 | 104.38 | 119.88 |
| **Current smoking** | binary | 414793 | no | no | no | no | yes |
| **Past smoking** | binary | 412557 | no | no | yes | yes | yes |
| **Family history of CAD** | categorical | 359472 | no | no | no | yes (1)[a] | yes (2)[b] |
| **Family history of DM2** | categorical | 385973 | no | no | no | no | mother |
| **Family history of high blood pressure** | categorical | 389301 | no | no | no | father | father and mother |
| **Family history of stroke** | categorical | 386630 | no | no | no | father | mother |
| **Physical activity (MET x hours/week)** | continuous | 379178 | 5.78 | 16 | 32 | 63 | 175.1 |
| **Coffee consumption (cups)** | continuous | 415021 | 0 | 0 | 2 | 3 | 6 |
| **Alcohol score** | continuous | 288169 | 0 | 0 | 2.5 | 10 | 10 |
| **AHEI score** | continuous | 199435 | 2.5 | 10 | 10 | 20.08 | 44.5 |
| **Surgery history** | binary | 322522 | no | no | no | no | yes |
| **Hormone replacement therapy** | categorical | 403518 | no | no | no | no | recent user (<3 years) |
| **Hypercholesterolemia medication excluding aspirin** | binary | 416936 | no | no | no | no | yes |
| **Sleep apnea** | binary | 416936 | no | no | no | no | no |
| **Irritable bowel syndrome** | binary | 416936 | no | no | no | no | no |
| **Heart valve problem** | binary | 416936 | no | no | no | no | no |
| **Arrhythmia** | binary | 416936 | no | no | no | no | no |
| **Congestive heart failure** | binary | 416936 | no | no | no | no | no |
| **Hyperthyroidism** | binary | 416936 | no | no | no | no | no |
| **Education** | categorical | 408500 | no | professional | professional | college or university | college or university |
| **Income (£)** | categorical | 353335 | <18,000 | 18,000 - 30,999 | 31,000 - 51,999 | 52,000 - 100,000 | >100,000 |
| **Insomnia** | categorical | 415605 | never/rarely | sometimes | sometimes | usually | usually |
| **Sleep duration (hours)** | categorical | 416117 | >4 and <6 or >9 and <11 | >=6 and <7 or >8 and <=9 | >=7 and <=8 | >=7 and <=9 | >=7 and <=10 |
| **Lymphocyte** | categorical | 395894 | >0.8 and <4.8 | >0.8 and <4.8 | >0.8 and <4.10 | >0.8 and <4.10 | >0.8 and <4.10 |
| **Monocyte** | categorical | 395894 | >0.2 and <0.9 | >0.2 and <0.9 | >0.2 and <0.9 | >0.2 and <0.9 | <=0.2 |
| **MCH** | categorical | 396632 | >=27 and <=34 | >=27 and <=34 | >=27 and <=34 | >=27 and <=34 | >34 |
| **Platelet** | categorical | 396631 | >=150 and <= 440 | >=150 and <= 440 | >=150 and <= 440 | >=150 and <= 440 | >=150 and <= 440 |
| **RDW** | categorical | 396633 | >= 11.6 and <= 14.6 | >= 11.6 and <= 14.6 | >= 11.6 and <= 14.6 | >= 11.6 and <= 14.6 | >14.6 |
| **CAD age** | continuous | 13607 | 45 | 54 | 58 | 62 | 67 |
| **DM2 age** | continuous | 8316 | 43 | 53 | 58 | 63 | 67 |
| **HPT age** | continuous | 36546 | 44 | 54 | 59 | 63 | 67 |
| **DVT age** | continuous | 7379 | 41 | 51 | 58 | 62 | 67 |

8

119    Types of variables, number of UKBB participants for each variable, and mean values (mode

120    categories for categorical and binary variables) for different percentiles are shown. The

121    number of participants for the CMD age variables corresponds to the number of prevalent

122    cases.

123    [a]Either father, mother, or sibling

124    [b]Any combination of two of the following: father, mother, or sibling

125

126        Physical activity was assessed as the metabolic equivalent of task (MET) calculated in

127    hours/week according to the "Guidelines for Data Processing and Analysis of the International

128    Physical Activity Questionnaire (IPAQ) [9]. MET coefficients are indicated in Table 1.

129    Alcohol score was calculated according to Alternative Healthy Eating Index (AHEI)

130    guidelines [10]. One alcohol serving corresponded to 11.4 grams of alcohol. Further, a

131    nutrition AHEI score was calculated as a sum of scores for the following nutrition categories:

132    vegetables, fruits, grains, sugar sweetened beverages and fruit juices, nuts, meat, fish, PUFA,

133    and alcohol. The nutrition scores were calculated according to AHEI guidelines [10].

134        In addition to the predicted CMD (target CMD), participants could of course experience

135    other competing CMD outcomes. We used the age of experiencing these non-target diseases

136    as an additional risk factor. For participants that did not experience a CMD event before

137    baseline (CMD-negative cases), the age of CMD was set to 100. This approach allowed for

138    incorporating time-dependent data without using the limitations of a modification of the Cox

139    model, such as a Cox proportional hazards time varying model, which is often used to address

140    time-dependency of predictors.

## Imputation of missing values

Multiple imputation by chained equations (MICE) implemented in Python (fancyimpute 0.3.1) and Bayesian ridge regression with the regularization parameter lambda of 0.001 was used for the imputation of missing values of continuous variables [11]. Parameters included initial filling with mean values, monotone visit sequence, the number of imputations = 100, the number of burn-in iterations = 10, no maximum and minimum possible imputed values, imputing with samples from posterior predictive distribution, the number of nearest neighbors for probabilistic moment matching = 5, and use of all columns to estimate current column. Cases with missing values in categorical variables were dropped before the imputation, and continuous variables were scaled to a range between 0 and 1.

## Variable selection for predictive modeling

Several approaches were employed for selecting variables included in the prediction model. Multicollinearity was first identified using pairwise correlation matrix (pandas 0.20.1), and the variables with the Pearson correlation coefficient higher than 0.3 were removed from the dataset. Recursive variable elimination with stratified 2-fold cross-validation (RFECV) on training datasets was then used to determine optimal number of variables by recursively considering smaller and smaller sets of variables (scikit-learn 0.20.0). One variable was removed at each iteration, minimum number of variables to be selected was one, and accuracy was used for scoring.

RFECV was used in combination with balanced random forest (imbalanced-learn 0.4.2) bivariate classification model. Parameters of the random forest model included the number of estimators = 100, Gini impurity as the quality of split, 'auto' sampling strategy, maximum

163    depth of the decision tree = 0, minimum number of samples required to split an internal node

164    = 2, minimum number of samples required to be at a leaf node = 1, minimum weighted

165    fraction of the sum total of weights required to be at a leaf node = 0, the number of variables

166    to consider when looking for the best split = 'auto', unlimited number of leaf nodes, minimum

167    impurity decrease threshold for node splitting = 0, bootstrapping, random sampling without

168    replacement, no use out-of-bag samples to estimate the generalization accuracy, the number

169    of jobs to run in parallel for both fit and predict = 1, resampling all classes, but the minority

170    class, the verbosity of the tree building process = 0, and balanced class weights.

171    In addition, principal component analysis (PCA) was used to validate the selection of

172    variables and to avoid overfitting and poor calibration by determining that the number of

173    selected variables is similar to the optimal number of principal components (scikit learn

174    0.20.0). The number of components to be retained was determined by using maximum-

175    likelihood density estimation and full singular value decomposition (utilizing LAPACK

176    library solver) as parameters of the PCA function, which applies Bayesian model selection to

177    probabilistic PCA in this configuration [12].

## Predictive models and performance metrics

179    Linear Cox proportional hazard (PH) models and non-linear ensemble survival models

180    were developed using lifelines 0.13.0 and scikit-survival 0.5 Python libraries, respectively.

181    Two types of non-linear models were developed: decision tree-based gradient-boosting using

182    Cox PH loss and gradient boosting with component-wise cubic smoothing splines as base

183    learners.

11

184    Discriminative ability of the risk prediction models was assessed by Harrell's

185    concordance index (c-index) [13, 14, 15] calculated for testing datasets as the proportion of all

186    comparable pairs in which the predictions and outcomes were concordant. Case pairs were

187    comparable if at least one of them was CMD-positive. If the estimated risk was larger for the

188    case with a lower time of event/censoring, the prediction of that pair was counted as

189    concordant. If predictions were identical for a pair, 0.5 was added to the count of

190    concordance. A pair was not comparable if an event occurred for both of them at the same

191    time or an event occurred for one of them, but the time of censoring was smaller than the time

192    of event of the first one. Prognostic indexes were used for the calculation of c-index.

193    In addition to c-index, we also used an additional metric for assessing the discriminative

194    ability of Cox PH models, which was based on statistical 'distance' between the probabilities

195    of experiencing a CMD event at certain time predicted for individuals from CMD-positive

196    and CMD-negative groups. In the 'distance' approach, statistical significance of the difference

197    between the two groups of probabilities was determined using one-way ANOVA. The result

198    of this test was reported as an *F*-statistic with corresponding *p*-value.

199    Calibration of Cox PH models was evaluated by the Hosmer-Lemeshov goodness-of-fit

200    test [16] and a calibration plot. The Hosmer-Lemeshow test was computed by partitioning the

201    testing set into decile groups based on the predicted absolute risk of CMD events at time

202    horizon of 5 years. Then, the number of CMD-positive and CMD-negative cases and the sum

203    of the predicted probabilities for the both types of cases was calculated in each group as

204    observed and not observed, and expected and not expected numbers, correspondingly. The

205    Hosmer-Lemeshow test statistic was calculated using the following formula:

206
$$H = \sum_{q=1}^{10} \frac{(Observed.A - Expected.A)^2}{Expected.A} + \frac{(Observed.not.A - Expected.not.A)^2}{Expected.not.A}$$

207    The resulted chi-square statistic was assessed using 8 degrees of freedom and was reported

208    with $p$-value. A calibration plot was created by plotting the predicted risk probabilities against

209    the observed risks for each group.

# Results

210

211    The study characteristics and the prevalence of six CMD at baseline for 416,936 UKB

212    participants that include CMD-positive cases that were identified by ICD-10 codes, self-

213    reports, or medication and had the date of the event determined based on the HES data are

214    shown in Tables 1 and 2. Average age of men and women in this population was $56.3 \pm 8.3$

215    and $56 \pm 8.1$ years, correspondingly. During follow-up (median 6.1 years), 98,254 incident

216    CMD events occurred in 67,785 participants that were free from the disease at baseline (Table

217    2).

218    **Table 2. Prevalent and incident events for various CMD.**

|  | Men | | Women | |
|---|---|---|---|---|
|  | **Prevalent events** | **Incident events** | **Prevalent events** | **Incident events** |
| **CAD** | 9442 (5.11%) | 9560 (5.17%) | 4165 (1.79%) | 5479 (2.36%) |
| **HTN** | 19489 (10.54%) | 27939 (15.11%) | 17057 (7.35%) | 24724 (10.66%) |
| **DM2** | 5155 (2.79%) | 7590 (4.1%) | 3161 (1.36%) | 5209 (2.25%) |
| **Stroke** | 740 (0.4%) | 1866 (1.01%) | 446 (0.19%) | 1290 (0.56%) |
| **DVT** | 3870 (2.09%) | 7387 (4.0%) | 3509 (1.51%) | 6447 (2.78%) |
| **AAA** | 241 (0.13%) | 644 (0.35%) | 38 (0.016%) | 119 (0.051%) |

219    The prevalence of CMD at the baseline and incidence of CMD during the follow-up are shown in

220    parenthesis.

221

# Imputation of missing data

222

14

223     Initial data quality evaluation showed that the number of missing values for examined

224     variables (Table 1) varied from 0 to ~52% with the mean of 6.3%, resulting in the no-null

225     values dataset sizes of ~78K – 81K (vs. initial ~380K – 416K).  As discussed in the methods,

226     imputation of missing values for all continuous variables (Table 1) excluding CMD age

227     variables, increased the sizes of CMD-specific datasets for predictive modeling to up to

228     ~195K – 215K.  The discriminative ability of the CAD risk model trained on the imputed

229     dataset with the sample size of 165,877 was tested on both imputed and unimputed datasets

230     with the same sample size of 41,470 to validate the imputation.  C-indexes calculated on the

231     imputed and unimputed testing sets were 0.787 and 0.803, implying higher discriminative

232     ability of the CAD model when tested on original, unimputed data.

## Predictive modeling

233

234     The discriminative ability of all Cox PH CMD models trained on the general population

235     after the imputation of missing data varied between the diseases with highest and lowest c-

236     indexes of 0.88 and 0.748 for AAA and DVT, respectively (Table 3).  Cox PH models were

237     further applied to calculate the risk probabilities of occurrence of a CMD event at 5 years

238     following the initial observation.   This time-to-event prediction was evaluated through

239     determination of the statistical 'distance' between CMD-positive and CMD-negative test

240     subgroups' risk scores (Table 3).  $F$-statistic values for the CMD models were highest for the

241     models with high discriminative ability, except for the AAA model due to the low prevalence

242     of this disease.

243

244

245 **Table 3: Performance of CMD risk prediction models.**

|  | C-index | Hosmer-Lemeshov test | | ANOVA test | |
|---|---|---|---|---|---|
|  |  | chi-2 | *p*-value | *F*-statistic | *p*-value |
| **CAD** | 0.787 | 55 | < 0.0001 | 24.7 | 1.80E-04 |
| **HPT** | 0.817 | 155 | < 0.0001 | 44.6 | 8.04E-07 |
| **DM2** | 0.873 | 54 | < 0.0001 | 36.6 | 1.20E-06 |
| **Stroke** | 0.783 | 18 | 0.02 | 17.6 | 6.20E-03 |
| **DVT** | 0.748 | 45 | < 0.0001 | 18.7 | 5.00E-03 |
| **AAA** | 0.88 | 17 | 0.03 | 15 | 1.20E-03 |

246 Performance is by c-index (discrimination), Hosmer-Lemeshov test (calibration), and the

247 statistical 'distance' approach based on one-way ANOVA test (discrimination of risk

248 probabilities). CMD-positive and negative groups were bootstrap sampled with replacement

249 (N=100) to provide comparable *F*-statistic (*p*-values) across different disease endpoints.

250

251      Probability density function, which specifies the probability of predictions falling within

252 a particular range of values for individuals from CMD-positive and CMD-negative test

253 subgroups (Fig 1) was used for the visualization of the statistical 'distance' approach. The

254 probability density function of the risk scores, as well as their distributions derived from

255 different CMD models demonstrated that the range of risk scores for the CMD-positive

256 subgroup was higher than that for the CMD-negative subgroup, and increased for CMD

257 models characterized by higher c-index. Higher ratio between maximum values of the two

258 probability density functions corresponded to higher discriminative ability.
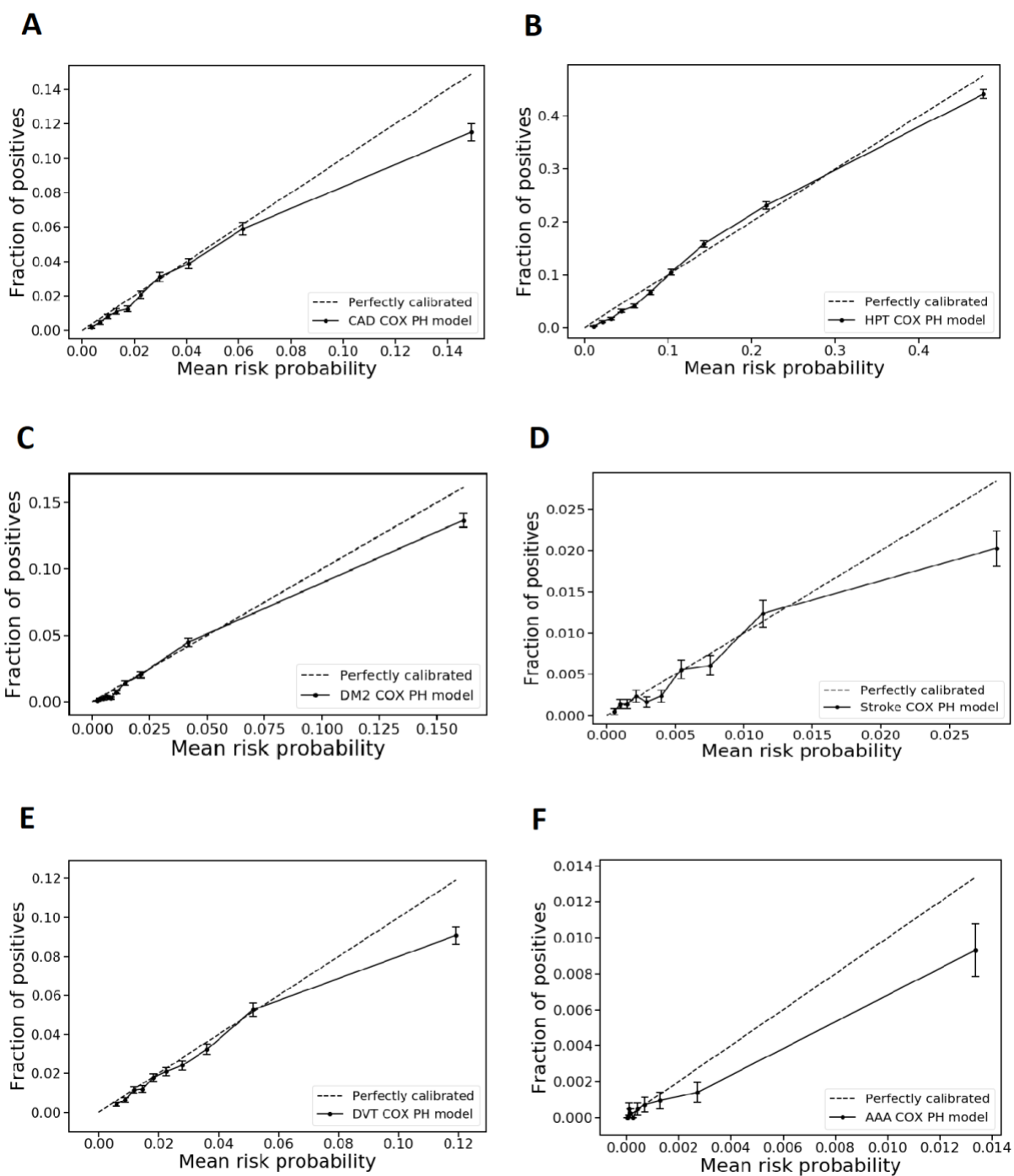
# Figure 1



259

260     **Fig 1. Statistical 'distance' approach.** Probability density function expressed in relation to

261     risk scores for six diseases (A-F) comparing participants developing CMD (CMD-positive,

262     _1) and those who did not develop (CMD-negative, _0) within 5 years of follow up.

263

264     Assessment of the calibration properties for the CMD predictive models as calculated by

265     the Hosmer-Lemeshow test (Table 3) and visualized by the calibration plot (Fig 2) showed

266     adequate overall calibration, but moderate overestimation of CMD risk in the highest decile of

267     risk scores.

# Figure 2



268

269    **Fig 2. Calibration plots for CMD prediction models.** Risk probabilities for six diseases (A-

270    F), were split into deciles and mean risk probability for each decile was plotted vs. the portion

271    of positive CMD cases in the decile for time horizon of 5 years.

272

273    In this study, the predictive performance of linear Cox PH models was compared with

274    ensemble non-linear models as discussed in the methods. Non-linear survival models

275    demonstrated comparable performance with the linear Cox model; however, this required

276    significantly more computation time.

## CMD risk factors

278    To better understand the contribution of various risk factors to the pathophysiology of

279    CMD, we ranked predictors of the risk of various CMD by the values of their regression

280    coefficients (Table 4), indicating the degree of the association between the predictor and the

281    outcome. Predictors presented in Table 4 represented only those with absolute values of

282    coefficients larger than 0.8 and *p*-values less than 0.001 (see S1 Table for all coefficients).

283    Statistical significance depended on the sample size and was affected by the prevalence of

284    CMD. Accordingly, the number of predictors varied for each disease model.

285

286

287

288

289

290 **Table 4. Ranked regression coefficients of predictors of the risk of various CMD models.**

|  | Variable | Coefficient | lower 95% CI | upper 95% CI | *p-value* |
|---|---|---|---|---|---|
| **CAD** | Forced expiratory volume | -3.45 | -4.08 | -2.82 | 3.18E-27 |
|  | Body mass index | 2.94 | 2.62 | 3.25 | 3.69E-76 |
|  | Age | 2.29 | 2.15 | 2.43 | 8.19E-225 |
|  | Heart valve problem | 0.99 | 0.82 | 1.16 | 3.15E-29 |
|  | Sex | 0.94 | 0.88 | 1.01 | 3.15E-156 |
|  | Family history of CAD (both parents) | 0.87 | 0.73 | 1.00 | 4.50E-36 |
|  | Hypercholesterol medication | 0.84 | 0.78 | 0.89 | 9.91E-189 |
| **HPT** | Diastolic blood pressure | 4.72 | 4.58 | 4.87 | 0.00E+00 |
|  | Body mass index | 3.69 | 3.54 | 3.83 | 0.00E+00 |
|  | Forced expiratory volume | -3.33 | -3.67 | -2.98 | 3.47E-81 |
|  | Age | 2.43 | 2.35 | 2.51 | 0.00E+00 |
|  | Coffee consumption | -1.61 | -2.11 | -1.11 | 2.72E-10 |
|  | Congestive heart failure | 1.32 | 0.93 | 1.70 | 2.45E-11 |
|  | Hypercholesterol medication | 1.20 | 1.17 | 1.23 | 0.00E+00 |
|  | CAD age | -1.16 | -1.25 | -1.06 | 6.57E-124 |
| **DM2** | Body mass index | 6.99 | 6.75 | 7.23 | 0.00E+00 |
|  | Forced expiratory volume | -6.54 | -7.23 | -5.85 | 2.11E-77 |
|  | MET hours | -1.84 | -2.76 | -0.92 | 9.17E-05 |
|  | Hypercholesterol medication | 1.82 | 1.76 | 1.89 | 0.00E+00 |
|  | Coffee consumption | -1.66 | -2.64 | -0.68 | 9.27E-04 |
|  | Age | 1.45 | 1.30 | 1.61 | 1.55E-77 |
|  | Family history of DM2 (both parents) | 1.40 | 1.26 | 1.54 | 2.04E-85 |
|  | AHEI score | 0.93 | 0.69 | 1.16 | 4.87E-15 |
| **Stroke** | Forced expiratory volume | -5.34 | -6.73 | -3.96 | 3.93E-14 |
|  | Age | 3.17 | 2.83 | 3.50 | 7.46E-77 |
|  | Diastolic blood pressure | 2.26 | 1.67 | 2.86 | 1.21E-13 |
|  | DVT age | -1.14 | -1.50 | -0.77 | 1.25E-09 |
|  | Diabetes age | -0.87 | -1.22 | -0.53 | 8.06E-07 |
|  | AHEI score | 0.85 | 0.37 | 1.33 | 5.58E-04 |
| **DVT** | Forced expiratory volume | -3.55 | -4.20 | -2.89 | 2.00E-26 |
|  | Body mass index | 2.58 | 2.26 | 2.90 | 3.42E-57 |
|  | Age | 1.94 | 1.80 | 2.08 | 3.18E-156 |
| **AAA** | Forced expiratory volume | -5.99 | -8.64 | -3.33 | 9.78E-06 |
|  | Age | 5.20 | 4.43 | 5.97 | 3.54E-40 |
|  | AHEI score | 1.98 | 1.12 | 2.83 | 5.83E-06 |
|  | Heart valve problem | 1.52 | 0.99 | 2.04 | 1.67E-08 |
|  | Sex | 1.47 | 1.06 | 1.88 | 1.98E-12 |
|  | Current smoking | 1.15 | 0.89 | 1.40 | 3.68E-18 |
|  | Hypercholesterol medication | 0.90 | 0.65 | 1.15 | 1.27E-12 |

291 Positive and negative signs indicate that corresponding factors increase or decrease the risk of

292 CMD, respectively. For the purpose of better presentation, only coefficients with absolute

293 values larger than 0.8 and *p*-values less than 0.001 are presented.

21

294         Across all disease models, age and low forced expiratory volume (FEV1) ranked as the

295         most important predictors. Higher body mass index (BMI) and hypercholesterolemia

296         medication were also among the strongest predictors for several models. Sex was ranked high

297         only for the CAD and AAA, which is in a good agreement with our observation that the

298         prevalence of these diseases was higher in men than in women. Family history ranked high

299         only in predicting CAD and DM2. Nutrition was among the most important predictors for

300         DM2, stroke, and AAA, which is likely explained by a healthier diet among individuals with

301         certain risk factors and predispositions. Similarly, coffee consumption was an important

302         predictor of HTN and DM2, possibly due to lower consumption in individuals with specific

303         risk factor profiles. Physical activity was an important predictor only for DM2, and younger

304         age of first occurrence of CAD, DVT and DM2 was among most important predictors for

305         HTN and stroke, respectively.

306      ## Validation

307         C-indexes for corresponding risk prediction benchmark models, with age and sex as the

308         only predictors, were lower (delta, 0.04 – 0.2) when compared to those of our newly

309         developed models. Broad range applicability and consistency of the performance of the

310         developed risk prediction models for each disease were further determined by assessing the

311         discriminative ability across subpopulations (Table 5). These subpopulations included (1)

312         'healthy' participants without any of the six target CMD at the baseline; (2) participants with

313         at least one pre-existing non-target CMD at the baseline; and (3) various age categories. The

314         performance of the models was highest in younger age and the healthy subgroup; while it

315         significantly dropped in the subpopulation with pre-existing CMD.

316 **Table 5. Validation of CMD risk prediction models.**

| Subpopulation | CAD | | HPT | | DM2 | | Stroke | | DVT | | AAA | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Cases, % | C-index | Cases, % | C-index | Cases, % | C-index | Cases, % | C-index | Cases, % | C-index | Cases, % | C-index |
| General (benchmark model) | 3.6 | 0.716 | 13.8 | 0.689 | 3 | 0.673 | 0.7 | 0.712 | 3.5 | 0.678 | 0.19 | 0.837 |
| Healthy + target CMD | 3.2 | 0.785 | 12.9 | 0.813 | 2.4 | 0.883 | 0.6 | 0.772 | 2.8 | 0.722 | 0.12 | 0.874 |
| Unhealthy + target CMD | 9.3 | 0.656 | 43.5 | 0.693 | 7.9 | 0.724 | 1.8 | 0.677 | 9.5 | 0.62 | 0.66 | 0.794 |
| CAD | 100 | n/a | 56.8 | 0.642 | 11.2 | 0.697 | 2.5 | 0.684 | 12.1 | 0.568 | 1.22 | 0.817 |
| HPT | 9.4 | 0.65 | 100 | n/a | 7.8 | 0.72 | 1.7 | 0.655 | 9 | 0.637 | 0.74 | 0.775 |
| DM2 | 14.5 | 0.631 | 52.5 | 0.624 | 100 | n/a | 3.3 | 0.662 | 11.3 | 0.57 | 0.34 | 0.812 |
| DVT | 8.8 | 0.695 | 26.1 | 0.734 | 7.5 | 0.752 | 2.3 | 0.733 | 100 | n/a | 0.7 | 0.907 |
| Age < 45 | 0.9 | 0.842 | 3.3 | 0.864 | 0.8 | 0.872 | 0.13 | 0.676 | 0.9 | 0.669 | 0.04 | 0.872 |
| Age 45-55 | 1.9 | 0.769 | 7.6 | 0.824 | 1.9 | 0.894 | 0.4 | 0.744 | 2 | 0.711 | 0.02 | 0.725 |
| Age 55-65 | 4.4 | 0.736 | 16.9 | 0.774 | 3.4 | 0.85 | 0.7 | 0.744 | 3.8 | 0.704 | 0.19 | 0.843 |
| Age 65-75 | 7.7 | 0.707 | 26.5 | 0.736 | 5.2 | 0.825 | 1.9 | 0.661 | 7 | 0.665 | 0.54 | 0.823 |

317 The performance of CMD models was tested on four different age group subpopulations.

318 Healthy subpopulation included individuals without *any* CMD at the baseline. Unhealthy

319 subpopulation included cases with any non-target CMD at the baseline.

320

## Discussion

### Principal findings

In this study, development and validation of a risk assessment platform applicable to six CMD is presented. The population-specific modeling for this platform was done using a dataset from the UK Biobank – a very large, longitudinal cohort study. This allowed us to derive prediction models and identify the most important contributing risk factors even for diseases with low incidence. Inclusion of a broad spectrum of risk factors allowed for modification of the array of input variables for the CMD risk prediction models included into the platform without significant decrease in their predictive performance. The models performed with high discriminative ability as demonstrated through extensive validation for different disease and age group subpopulations. Accordingly, this platform can accommodate different types of data sets and is applicable to population analysis, as well as individual assessment.

There is an abundance of risk predictors for CMD, and multiple prior attempts of combining them into risk calculators [17-19]. One of the major impediments for wide-spread application of these risk predictors includes lack of uniform validation through large population analyses. A comprehensive review found 363 models for cardiovascular risk stratification that have been developed and reported [20]. Only a minor collection of these models had sufficient evaluation according to contemporaneous analysis standards for either development or validation. For example, 39% of the 363 models analyzed utilized C-statistics for their development, and just over 60% for their validation. An even smaller number of the models utilized calibration as any part the performance measures. Although, the more recent

24

343     models (since 2009) were more consistent in providing performance reports: 76% as part of

344     their development, and up to 90% as part of validation [20].

345     In the current study, the discriminative ability of the developed models was similar or

346     exceeded established models when available.  For example, the Framingham Risk Score for

347     coronary artery disease have been determined to be close to 0.76 and 0.79 for men and women,

348     respectively [21]; these reported results were obtained only in the presence of all of the

349     laboratory data and for a pre-selected small population.  The modeling described for the platform

350     in this report allows for incorporation of contemporary risk information.  This is becoming

351     increasingly important, since such more limited risk calculators may fail to express the accurate

352     and true risk for a significant population.  As demonstrated previously, either 50% of patients

353     with CMD lack conventional risk factors or the conventional risk factors fail to explain more

354     than 15-50% of the incidence of CHD [22-26].

355     The ability to incorporate socioeconomical data and nutritional information collectively

356     can complement the basic information that is equivalent to conventional biomarkers.  This is

357     demonstrated in this study, as the performance of the current platform was achieved without

358     the utilization of the blood laboratory information, such as lipid levels or blood glucose levels

359     (as those were not available in UKBB at the time of this study).  Utilization of a polygenic

360     scoring is underway and can reveal a population at risk or protected from development of

361     CMD [27-29].  It is expected that incorporation of the polygenic scoring will further increase

362     the predicative performance of the current platform.

363     **Limitations of this study**

364     Considering the fact that the UKBB population is not a complete representative of the

365     UK or US populations, the main limitation of this study is that the developed models may

25

366  need to be examined with inclusion of more diverse population. Predictive performance of

367  the models was higher when tested on healthier and younger subpopulations. At the same

368  time, training and calibration on CMD-specific datasets are required to improve

369  discriminative ability of the models across CMD subpopulations. Considering the fact that

370  the datasets used in predictive modeling were almost identical for different CMD, various

371  predictive performances of the CMD models imply that despite overlapping

372  pathophysiological pathways for various CMD, there are predictors specific for different

373  CMD.

374  **Future directions**

375  Considering computational limitations of non-linear survival models, bivariate time-

376  dependent classification models utilizing machine learning algorithms can be used in future

377  for determining the probability of CMD events at certain time horizons. The availability of

378  relatively large healthcare datasets will further support the application of deep learning in

379  time-dependent risk predictive modeling feasible. Incorporation of genetic and other -omics

380  data may further improve the predictive functionality provided by this platform.

381  # Conclusions

382  In this report, we present development and validation of a new generation of disease risk

383  prediction models. The differentiation variables of this platform include: a) assessment of

384  multiple related diseases according to their associated outcomes (not just coronary artery

385  disease); b) inclusion of contemporary risk factors; c) variable engineering and processing

386  that allows for inclusion of data from different sources and addressing missing data points; d)

387  population-specific stratification to assess risk prediction in different subgroups; e) being

388  modular in nature to allow for inclusion of other risk determinants, such as genetic

26

389    information; and f) being applicable at individual, as well as population level. These

390    variables were designed into the platform in order to provide applicability of risk prediction to

391    managing and changing the course of cardiometabolic diseases.

392

393

394    **Acknowledgements**

# References

1. Benjamin EJ, Virani SS, Callaway CW, Chang AR, Cheng S, Chiuve SE, et al. The American Heart Association Council on Epidemiology and Prevention Statistics Committee and Stroke Statistics Subcommittee. H2018 update: a report from the American Heart. Circulation 2018 Jan 31. DOI: 10.1161/CIR.0000000000000558.

2. Muka T, Imo D, Jaspers L, Copani V, Chaker L, vad der Lee SJ, et al. The global impact of non-communicable diseases on healthcare spending and national income: a systematic review.Eur J Epidemiol. 2015 30(4): 251-77.

3. Neumann PJ, Cohen JT, Cost savings and cost-effectiveness of clinical preventive care. Synth Proj Res Synth Rep. 2009 (18). pii: 48508. doi: 48508. Epub 2009 Sep 1.

4. Anderson KM, Odell PM, Wilson PW, Kannel WB. Cardiovascular disease risk profiles. Am Heart J. 1991; 121: 293-8. doi:10.1016/0002-8703(91)90861-B. pmid:1985385.

5. Conroy RM, Pyörälä K, Fitzgerald AP, et al. SCORE project group. Estimation of ten-year risk of fatal cardiovascular disease in Europe: the SCORE project. Eur Heart J. 2003; 24:987-1003.

6. Hippisley-Cox J, Coupland C, Robson J, Brindle P. Derivation, validation, and evaluation of a new QRISK model to estimate lifetime risk of cardiovascular disease: cohort study using QResearch database. BMJ 2010;341:c6624. doi:10.1136/bmj.c6624. pmid:21148212.

7. Wang J, Tan G, Han L, Bai Y, He M, Liu H. Novel Biomarkers for Cardiovascular risk prediction. J Geriatr Cardiol. 2017 Feb; 14(2): 135–150.

8. Palmer LJ. UK Biobank: bank on it. Lancet 2007; 369: 1980-1982.

417    9.  Craig CL, Marshall AL, Sjöström M, Bauman AE, Booth ML, Ainsworth BE, et al. International

418        physical activity questionnaire: 12-country reliability and validity. Med Sci Sports Exerc.

419        2003; 35(8):1381-95.

420    10. Chiuve SE, Fung TT, Rimm EB, Hu FB, McCullough ML, Wang M, et al. Alternative dietary

421        indices both strongly predict risk of chronic disease. J Nutr. 2012; 142(6):1009-18.

422    11. White IR, Royston P, Wood AM. Multiple imputation using chained equations: issues and

423        guidance for practice. Stat Med. 2011; 30: 377-399.

424    12. Minka T. Automatic choice of dimensionality for PCA. No. 5141-16. M.I.T. media laboratory

425        perceptual computing section technical report; 2000.

426    13. Harrell FE, Califf RM, Pryor DB, Lee KL, Rosati RA. Evaluating the yield of medical tests.

427        Journal of the American Medical Association. 1982; 247:2543–46.

428    14. Harrell FE, Lee KL, Califf RM, Pryor DB, Lee KL, Rosati RA. Regression modeling strategies for

429        improved prognostic prediction. Statistics in Medicine. 1984; 3:143–52.

430    15. Harrell FE, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models,

431        evaluating assumptions and adequacy, and measuring and reducing errors. Statistics in

432        Medicine. 1996; 15:361–87.

433    16. Hosmer DW, Hosmer T, Le Cessie S, Lemeshow S. A comparison of goodness-of-fit tests for

434        the logistic regression model. Statistics in Medicine. 1997; 16:965-980.

435    17. Piepoli MF, Hoes AW, Agewall S, Albus C, Brotons C, Catapano AL, et al. 2016 European

436        Guidelines on cardiovascular disease prevention in clinical practice: The Sixth Joint Task

437        Force of the European Society of Cardiology and Other Societies on Cardiovascular Disease

438        Prevention in Clinical Practice (constituted by representatives of 10 societies and by invited

439     experts) Developed with the special contribution of the European Association for

440     Cardiovascular Prevention & Rehabilitation (EACPR). Atherosclerosis. 2016; 252:207–74.

441     18. Goff DC Jr., Lloyd-Jones DM, Bennett G, Coady S, D'Agostino RB Sr., Gibbons R, et al. 2013

442     ACC/AHA guideline on the assessment of cardiovascular risk: a report of the American

443     College of Cardiology/American Heart Association Task Force on Practice Guidelines. J Am

444     Coll Cardiol. 2014;63(25 Pt B):2935–59.

445     19. Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, Minhas R, Sheikh A, et al. Predicting

446     cardiovascular risk in England and Wales: prospective derivation and validation of

447     QRISK2. Bmj. 2008;336(7659):1475–82.

448     20. Damen JA, Hooft L, Schuit E, Debray TP, Collins GS, et al. Prediction models for

449     cardiovascular disease risk in the general population: systematic review. BMJ. 2016; 353:

450     i2416.

451     21. D'Agostino RB, Vasan RS; Pencina MJ, Wolf PA, Cobain M; Massaro JM et al. General

452     Cardiovascular Risk Profile for Use in Primary Care: the Framingham Heart Study.

453     Circulation. 2008; 117(6):743-53.

454     22. Hennekens CH. Increasing burden of cardiovascular disease: current knowledge and future

455     directions for research on risk factors.  Circulation.1998;97:1095-1102.

456     23. McKechnie RS, Rubenfire M. The role of inflammation and infection in coronary artery

457     disease: a clinical perspective.  ACC Curr J Rev. 2002;11:32-34.

458     24. Futterman LG, Lemberg L. Fifty percent of patients with coronary artery disease do not have

459     any of the conventional risk factors.  Am J Crit Care.1998;7:240-244.

460     25. Lefkowitz RJ, Willerson JT. Prospects for cardiovascular research.  JAMA.2001;285:581-587.

461    26. Khot UN, Khot MB, Bajzer CT, Sapp SK, Ohman M, Brener SJ, et al. Prevalence of

462        Conventional Risk Factors in Patients with Coronary Heart Disease JAMA. 2003;290(7):898-

463        904.

464    27. Kathiresan S, Melander O, Anevski D, Guiducci C, Burtt NP, Roos C, et al. Polymorphisms

465        Associated with Cholesterol and Risk of Cardiovascular Events. N Engl J Med 2008;

466        358:2299-2300

467    28. Khera AV, Emdin CA, Drake I, Natarajan P, Bick AG, et al. Genetic Risk, Adherence to a

468        Healthy Lifestyle, and Coronary Disease. N Engl J Med 2016; 375:2349-58.

469    29. Khera AV, Chaffin M, Aragam KG, Haas ME, Roselli C, Choi SH, et al. Genome-wide polygenic

470        scores for common diseases identify individuals with risk equivalent to monogenic

471        mutations. Nat Genet. 2018; 50(9):1219-1224.

472

473

# Supporting information

475    **S1 Table. Cox PH model regression coefficients for six CMD.** Regression coefficients (coef) and

476    corresponding standard errors (se), $p$-values, lower and upper 95% confidence intervals are

477    presented.

478

31

483    decision to publish, or preparation of the manuscript. MR did not receive any financial
484    compensation for participation. The specific roles of these authors are articulated in the
485    'Author contributions' section.
486

487    **Author Contributions:**

488        1.  Conceptualization: MR AT

489        2.  Data curation: IP AT

490        3.  Formal analysis: IP

491        4.  Funding acquisition: MR

492        5.  Investigation: MR AT EI

493        6.  Methodology: MR AT IP AG EI

494        7.  Project administration: MR AT

495        8.  Resources: MR

496        9.  Software: AT IP

497        10. Supervision: MR AT

498        11. Validation: MR EI AG

499        12. Visualization: IP

500        13. Writing – original draft: AT IP MR

501        14. Writing – review & editing: MR EI AG

502

# Figure 2

# Figure 1