

1 MSGene: Derivation and validation of a multistate model for lifetime risk of coronary artery
2 disease using genetic risk and the electronic health record

3
4 Sarah M. Urbut, MD, PhD^{1,2,3}, Ming Wai Yeung, MSc⁴, Shaan Khurshid, MD, MPH^{2,5,6}, So Mi
5 Jemma Cho^{1,2,3,7}, Art Schuermans, BS^{2,3,8}, Jakob German, MSc^{9,10}, Kodi Taraszka, PhD¹¹, Akl C.
6 Fahed, MD, MPH^{1,2,3}, Patrick Ellinor, MD, PhD^{1,2,5,6}, Ludovic Trinquart, PhD¹², Giovanni
7 Parmigiani, PhD^{11,13}, Alexander Gusev, PhD^{11,14}, Pradeep Natarajan, MD, MMSc*^{1,2,3}

- 8
9 1. Division of Cardiology, Department of Medicine, Massachusetts General Hospital,
10 Harvard Medical School, Boston, MA
11 2. Program in Medical and Population Genetics, Broad Institute of MIT and Harvard,
12 Cambridge, MA
13 3. Center for Genomic Medicine, Department of Medicine, Massachusetts General
14 Hospital, Boston, MA
15 4. University of Groningen, University Medical Center Groningen, Department of
16 Cardiology, 9700 RB Groningen, The Netherlands
17 5. Demoulas Center for Cardiac Arrhythmias, Massachusetts General Hospital, Boston, MA
18 6. Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA
19 7. Integrative Research Center for Cerebrovascular and Cardiovascular Diseases, Yonsei
20 University College of Medicine, Seoul, Republic of Korea
21 8. Faculty of Medicine, KU Leuven, Leuven, Belgium
22 9. Institute for Molecular Medicine Finland (FIMM), HiLIFE, University of Helsinki, Helsinki,
23 Finland
24 10. Eric and Wendy Schmidt Center, Broad Institute of MIT and Harvard, Cambridge, MA,
25 USA
26 11. Dana Farber Cancer Institute, Boston, MA
27 12. Tufts University Medical Center, Boston, MA
28 13. Harvard School of Public Health, Boston, MA
29 14. Department of Medicine, Harvard Medical School, Boston, MA

30
31 *Corresponding author

32
33 Word Count:

34 Abstract: 150 words

35 Main text: 3565/4000

36 Main Figures and tables: 6

37 Supplementary figures: 17

38

39 Address for correspondence:

40

41 Pradeep Natarajan, MD, MMSc

42 185 Cambridge Street, CPZN 3.184

43 Boston, MA 02114

44 Tel: 617-726-1843 | E-mail: pnatarajan@mgh.harvard.edu | Twitter: @pnatarajanmd

45 **Abstract**

46 Currently, coronary artery disease (CAD) is the leading cause of death among adults worldwide.
47 Accurate risk stratification can support optimal lifetime prevention. We designed a novel and
48 general multistate model (MSGene) to estimate age-specific transitions across 10
49 cardiometabolic states, dependent on clinical covariates and a CAD polygenic risk score.
50 MSGene supports decision making about CAD prevention related to any of these states. We
51 analyzed longitudinal data from 480,638 UK Biobank participants and compared predicted
52 lifetime risk with the 30-year Framingham risk score. MSGene improved discrimination (C-index
53 0.71 vs 0.66), age of high-risk detection (C-index 0.73 vs 0.52), and overall prediction (RMSE
54 1.1% vs 10.9%), with external validation. We also used MSGene to refine estimates of lifetime
55 absolute risk reduction from statin initiation. Our findings underscore the potential public health
56 value of our novel multistate model for accurate lifetime CAD risk estimation using clinical
57 factors and increasingly available genetics.

58

59 **Introduction**

60
61 Coronary artery disease (CAD), remains the leading cause of morbidity and mortality
62 worldwide.¹ Estimating an individual's risk of developing CAD over the lifetime is essential for
63 timely and effective prevention and intervention.²⁻⁵ Traditional risk prediction models, such as
64 the Pooled Cohort Equations (PCE) 10-year risk score, have guided clinical decisions and
65 preventive strategies; however, these models come with inherent limitations.⁶⁻⁸ A 30-year or 10-
66 year window provides only a fixed, albeit extended, snapshot of risk. It neither captures the
67 entirety of an individual's lifetime risk nor provides dynamic, age-specific insights beyond these
68 arbitrary periods. Most importantly, there is a growing need for models capable of both
69 recognizing undertreated younger patients while reducing over-estimation in older patients.^{7,9,10}

70 Current guidelines^{9,11,12} recommend the consideration of primordial risk factors in risk-
71 stratifying patients, and call for better methods of estimating lifetime risk. Recent evidence
72 suggests that lifetime risk assessment provides a more comprehensive picture of an individual's
73 propensity for developing CAD across time.^{13,14} Traditional factors in combination with genomic
74 risk can confer a disproportionately elevated risk for CAD in the long term.^{2,15-17} Focusing on
75 lifetime risk allows for more effective patient counseling, tailored preventive measures, and
76 earlier interventions that may delay or prevent the onset of CAD altogether.^{18,19}

77 Because of the multifactorial nature of CAD, there is an increasing need for continuously
78 updated, dynamic and individualized CAD risk predictions that span a patient's entire life.^{2,14,20}
79 Such risk prediction models could improve the identification of undertreated younger patients
80 while avoiding risk over-estimation in older patients.^{7,9,10} Understanding risk from this
81 perspective allows for more informed and timely interventions, potentially even before the
82 conventional risk windows are applicable.

83 Here, we introduce the MSGene model — a multistate model designed to predict the
84 lifetime risk of CAD, conditional on both time-fixed and time-dependent variables. Multistate
85 models allow for the estimation of the risk of an individual transitioning between health states²¹⁻
86 ²⁵ through flexible estimation of conditional probabilities by modeling the transitions between
87 states over time. By modeling the different health states simultaneously, they naturally account
88 for competing risks.

89 MSGene is capable of modeling the dynamic transitions from risk factor states to CAD
90 with age-specific coefficients. Critically, our approach differs from a traditional Markov-based
91 multistate model^{21,22} by extending our model to the time inhomogeneous case and allowing our

92 transitions to vary with age, and also from traditional Cox models by allowing for non-
93 proportional hazards.

94 In the current study, we develop and validate the MSGene model. We evaluate the
95 performance compared to the traditionally employed Framingham 30-year²⁶ and PCE 10-year^{5,6}
96 models. We then estimate the potential ability of MSGene to reduce CAD events by guiding
97 timely initiation of statin therapy and demonstrate the benefit of a multistate framework to
98 incorporate dynamic changes in treatment decisions for unique patient profiles.

99 **Results**

100 **Novel multistate model with time-dependent transitions**

101 We build a novel time-dependent multistate model in which age is the time scale.
102 For each age and current state (**Fig. 1**), we model the one-year probability of transition from
103 state to state as a logistic regression conditional on both time-fixed covariates (sex, CAD-PRS),
104 and time-dependent covariates (smoking, use of anti-hypertensives or statins) (**Methods**). This
105 methodology defines an inhomogenous Markov transition model which can be used to compute
106 the probability of reaching any state of interest during one's lifetime, among other quantities.
107 Here, to compare our model to existing tools we focus on CAD.

108 We use a limited set of covariates (**Methods**) as a result of the variable selection
109 described in **Supp. Table 1**. To improve estimation efficiency, we smooth each set of state to
110 state coefficients across ages using a flexible tricube distance weighted least square local
111 regression²⁷ with inverse variance weighting of raw estimate. This allows for the sharing of
112 information across ages in instances in which the number of individuals at a particular transition
113 may be small. We calculate risk under a statin-treated and statin-untreated strategies by
114 imputing trial-imputed relative risk reduction of statin use on each annual age-specific transition
115 (**Methods**). We develop this in the R programming language (4.3.0) and provide detailed code
116 and an interactive application for users.

117 **Baseline characteristics**

118 We considered 480,638 individuals: 260653 (54.2%) were female with 43,855 (11.1%) incident
119 coronary artery disease diagnoses (**Table 1**) with a median 29.9 years [22.4–35.1] years of
120 follow-up and median age of first observation in EHR 24.3 [IQR: 18.0, 37.1] after excluding
121 20,534 who lacked sufficient covariates or had CAD at baseline (**Fig. 2**). MSGene allows for
122 visualization of the proportional representation by risk factor at each age (**Fig. 1**): approximately
123 39.6% are ultimately diagnosed with hypertension, 23.6% with hyperlipidemia, and 9.9% with

124 Diabetes mellitus (1 or 2). Furthermore, 10.5% report currently smoking and 20.3% began
125 antihypertensive use during the course of our study; 46.1% also contributed to the general
126 practice cohort, and the distribution of risk factors was homogenous between subsets (**Supp.**
127 **Fig. 1**). We use 80% of our data as training and 20% as testing (**Fig. 2**) for internal cross-
128 validation and to optimize model fit. Accordingly, this divides our data into a training set for
129 model fitting using 384,510 samples and a testing data set of 79,117 unique individuals. We
130 report the lifetime risk remaining at any age as one minus the product of the complement of the
131 interval age and state-specific transition to CAD probabilities.

132 **Modeling transitions**

133 Using our multistate approach, MSGene, we describe the overall state distribution across the
134 lifespan in our cohort, normalizing to exclude censoring at each age (**Fig. 1**). At age 40 years,
135 94.4% of individuals are in the healthy category, with 4.1% in the hypertensive category and
136 0.3% with a diagnosis of CAD. By age 76 years, CAD state occupancy peaks at 12.5% of
137 uncensored individuals, and health is reduced to 27.6% of uncensored individuals. By age 80
138 years, 7.4% have died.

139 **Improved detection of early events when compared to 10-year risk**

140 When compared to the PCE, a 10% lifetime threshold using MSGene uniquely identifies 5315
141 (59.3%) cases versus 123 (1.3%) cases using the 10-year PCE (5% threshold) alone at age 40.
142 This reduces to <1% of cases at age 68 (vs 81% with PCE) (**Supp. Fig 2**). At age 40, MSGene
143 had substantially greater sensitivity for lifetime CAD events compared to PCE (event
144 reclassification 58.2%, 95% CI 58.1–58.3%), at the cost of moderate inappropriate up-
145 classification of lifetime non-events (non-event reclassification –37.3%, 95% CI 37.2–37.4%). At
146 age 70, MSGene had substantially greater specificity compared to PCE (non-event
147 reclassification 32.1%, 95% CI 31.9–32.1%), at the cost of some inappropriate down
148 classification of events (event reclassification –12.5%, 95% CI –12.4 to –12.6%). Overall,
149 reclassification was consistently favorable (median NRI 0.12) over 40 years of consideration.
150 Furthermore, 9.7% (95% CI 9.6–9.8%) of individuals in the top 20% of genetic risk are identified
151 to have greater than 10% MSGene predicted lifetime risk, while only 3.1% (95% CI 2.9–3.2%) of
152 those in the bottom 20% of genetic risk achieve this level of risk (**Supp. Fig 2**).

153 **Improved calibration when compared to 30-year risk score**

154 MSGene had improved results when compared to FRS30RC. We compared the average
155 predicted risk by sex and genomic risk strata with empirical overall incidence rates. In healthy

156 individuals, the RMSE of MSGene is 1.06% (1.04% males, 1.09% females, SEM 0.06) while
157 FRS30RC is 10.9% (12.1% males, 10.1% females, SEM 0.07, **Supp. Fig. 3**). FRS30RC
158 increases monotonically across the lifespan. When restricting the analysis to ages 40 and 50 for
159 whom 30 years of follow-up is available, the RMSE is 0.98% with MSGene when compared to
160 5.68% for FRS30RC. We further compute the RMSE starting from additional single-risk factor
161 phenotype states (hypertension, hyperlipidemia, and diabetes) across a grid of covariate
162 choices (**Supp. Table 1**).

163 **Dynamic effects of 10-year, 30-year and remaining lifetime risk**

164 MSGene allows for the estimation of survival curves for an individual starting from a given age,
165 and for updated remaining lifetime curves asked over a range of ages. We compute the
166 remaining lifetime risk when compared with FRS30RC, as recalibrated for our population.²⁸
167 First, we depict the predicted survival curve for individuals of six different genetic and sex strata
168 starting in health at age 40. Under this traditional analysis, CAD-free survival is projected to
169 decline monotonically as a function of sex and genetic risk to 96.8% (95% CI 96.78–96.82) for a
170 female in the lowest genetic strata and to 81.26% (95% CI 81.24–81.28) for a male in the
171 highest genetic strata. However, a remaining lifetime risk curve reveals opposite behavior: for
172 example, a high genetic-risk male has a 22.9% (95% CI 22.7–23.1%) risk without treatment at
173 age 40, but the same high-risk male has only a 10.21% (95% CI 10.20–10.22%) risk of
174 developing CAD if he remains CAD-free at age 70. This contradicts the 10-year risk prediction,
175 in which 10-year risk rises from 2.84% at age 40 to 10.21% at age 70 (**Fig. 3, Supp. Tables 2–**
176 **17**). We compare this to FRS30RC projections²⁶ and note that while remaining lifetime risk
177 declines with age, the extended fixed-window (FRS30RC) approach shows monotonically
178 increasing risk across genetic strata. In our cohort the FRS30RC risk for a high genetic-risk
179 male rises from 13.4% at age 40 to 33.0% (**Fig. 3**) at age 70 using the recalibrated measure.
180 When applying trial-estimated statin benefit via introducing a trial-estimated relative risk
181 reduction to each annual transition probability²⁹ (**Methods, Eqn. 2**) under MSGene lifetime
182 projections, predicted absolute risk under treatment for the same high-genetic-risk male at age
183 40 improves from 22.86% (95% CI 22.85–22.87%) to 18.70% (95% CI 18.69–18.71%) over the
184 40-year span. This is compared to a smaller decline from 10.21% (95% CI 10.19–10.22%) to
185 8.25% (95% CI 8.24–8.26%) at age 70.

186 **Dynamic prediction: Model assessment**

187 An updated lifetime prediction, conditional on a patient's current state, can be made per year,
188 using age-specific coefficients. We use these updated predictions as covariates in a time-
189 dependent Cox model to evaluate the performance of our model on predicting time to event
190 **(Methods)**. We first consider the age distribution at which an individual first exceeded a lifetime
191 risk threshold of 10% using MSGene or FRS30RC, or using a PCE-derived 10-year risk
192 threshold of >5%. Using MSGene to assess lifetime risk, 44.8% percent of individuals exceed
193 this threshold at age 40 while 38.9% never do. With FRS30RC, 44.1% exceed this threshold at
194 age 40, but virtually all (99.8%) exceed this threshold by age 80. Using the first age exceeded
195 under each model as a time-dependent predictor of CAD status, we find that MSGene improves
196 model concordance by 21% (C-index 0.73 vs 0.52, $p < 2 \times 10^{-16}$) and of the 10-year index by
197 17.4% (C-index 0.55, $p < 2 \times 10^{-16}$) **(Fig. 4a-d)**.

198 We then use the yearly time- and state-varying predictions as predictors in a time-
199 dependent Cox proportional hazard model in which one's score is recorded annually in non-
200 overlapping intervals and estimate the concordance of this model. The concordance of this time-
201 dependent model using dynamic MSGene predictions exceeds that of the updated FRS30RC
202 predictions by 0.71 vs 0.66, $p < 2 \times 10^{-16}$ **(Fig. 4e-g)**. We repeat these results using the subset
203 with general practice (GP) records alone for both training (80%) and testing (20%) and the
204 results hold for both the thresholding analysis (C-index 0.71 vs 0.53, $p < 2 \times 10^{-16}$) and
205 continuous time-dependent analysis (C-index 0.73 vs 0.67, $p < 2 \times 10^{-16}$, **Supp. Fig. 4-5**).

206 **Estimated benefit**

207 Our model incorporates the estimated benefit of a treatment strategy, assessed conditional on
208 starting age and risk status. Using a randomized clinical trial (RCT)-imputed annual risk
209 reduction of 20% for statins on statin-free individuals,^{30,31} we observe an inverse relationship
210 between predicted 10-year risk and expected benefit. An individual with the highest genetic risk
211 at age 40 has a predicted 10-year risk (4.2%, SD 0.01) roughly equivalent to the lowest genetic-
212 risk individual at age 70 (3.9%, SD 0.01), but an expected lifetime absolute risk reduction of 5%
213 (SD 0.01) at age 40 versus only 0.8% (SD 5×10^{-2}) at age 70 **(Fig. 5)**. When we consider the
214 distribution of all starting states, we see that the mean absolute risk reduction is the greatest for
215 younger individuals (4.6–7.2%; SD 0.01) across risk states at age 40, to a mean absolute risk
216 reduction of 0.3–3.5% (SD 0.01) at age 79.

217 **Improvement in discrimination over the cumulative horizon**

218 When considering only the presence or absence of disease over observed time without regard
219 to timing, the AUC–ROC of a model comparing the prediction of cumulative occurrence using
220 updated MSGene lifetime score shows greater performance than that of either FRS30 or
221 FRSRC early in the life course (**Supp. Fig. 6**) (0.69 vs. 0.65, $p < 2 \times 10^{-16}$ at age 40) and also
222 based on precision-recall (0.20 vs 0.16 at age 40, $p < 0.01$). Both metrics exceeded the
223 estimation of lifetime risk using genetics as a predictor alone. In general, when comparing
224 individuals captured by MSGene but not by FRS30RC, MSGene identified more women and
225 individuals at higher genetic risk. With time, these differences were more profound (**Supp. Fig.**
226 **7**).

227 **External validation**

228 We then performed external validation of MSGene in the FOS cohort, using first measurements
229 to ensure optimal follow-up duration. FOS is a community-based cohort recruited in 1971 with a
230 median 39 years of follow-up [IQR 38–40], median age of enrollment 35 years [IQR 28–44]
231 (**Supp. Fig. 8**). MSGene again had favorable discrimination (age 40: 0.75 [95% CI 0.69–0.82]
232 vs. 0.73 [95% CI 0.66–0.80]; age 55: 0.63 [95% CI 0.42–0.84] vs. 0.53 [95% CI 0.29–0.76]) and
233 calibration (RMSE 8.4% vs. 11.3%, $p < 2 \times 10^{-16}$) when compared to FRS30 (**Supp. Fig. 9**).

234 **Discussion**

235 This study introduces a novel method called MSGene, which aims to assess the risk of
236 developing CAD and other health states over the lifespan. Our dynamic lifetime risk predictions
237 improve considerably calibration and discrimination and improve the identification of younger
238 individuals at high risk without overestimating risk in older adults, compared to previous models.
239 Our projected benefit analysis shows large reduction in preventable CAD events if statin therapy
240 is guided by MSGene.

241 The technique utilizes generalized linear models (GLMs) to compute the transition
242 probabilities between different states (e.g., from a healthy state or risk factor to CAD, death, or
243 intermediate risk) for every age over the observed life span. The novelty derives from four
244 features: 1) the provision of unique age-dependent models via GLMs that allow the relationship
245 of each covariate on the outcome to vary freely with time; 2) the calculation of risk conditional on
246 time-dependent states; 3) the assessment of a multistate model via time-dependent Cox
247 modeling; and 4) the unique use of the UKB EHR as a comprehensive longitudinal data
248 resource. The study follows individuals from adulthood through their enrollment in the linked

249 health record. By incorporating age and time dependence, this method provides annual risk
250 estimates that include the entire lifespan.

251 Over a lifetime horizon, the dynamic change in risk makes accurate lifetime risk
252 estimations challenging.^{4,7,11} However, leveraging genetics **in addition to multi-state**
253 **modeling**, MSGene enhances lifetime risk predictions, effectively identifying individuals
254 previously deemed low-risk. MSGene enhances lifetime risk predictions, effectively identifying
255 individuals previously deemed low-risk. The model's age-dependent features, producing age-
256 sensitive coefficients, negate the need to rely on fixed parametric interactions between each
257 covariate and time, a prevalent limitation in traditional models.⁶ We show that using updated
258 estimates conditional on the dynamic state of an individual improves *time-to-event* prediction
259 overall.

260 Through the incorporation of treatment, we show that those individuals with the greatest
261 and least expected absolute risk reduction from statin therapy actually have a similar 10-year
262 risk. However, this short-term focus is what current clinical methods rely upon.⁷ Presented
263 effects are conservative as statin effects may magnify with duration and on CAD PRS
264 background.^{19,32–34}

265 Our approach facilitates accurate event prediction both for undercaptured young
266 individuals and also lower-risk older individuals who might otherwise be included in a fixed-
267 window approach that extends the time horizon: our median global net reclassification when
268 compared with a 10-year approach is 12.2% [IQR 5.5–18.6%] over 40 years. This in part
269 explains the improvement in overall time-dependent performance when incorporated into a time-
270 to-event framework. Using a time-dependent evaluation, the distribution of the first age at which
271 a lifetime threshold is exceeded demonstrates that MSGene optimally identifies at-risk
272 individuals without indiscriminately calling all patients ‘at-risk’. However, future work is
273 warranted to determine optimal thresholds of lifetime risk to maximize potential benefits among
274 high-risk younger individuals while reducing unnecessary costs and harms to low-risk older
275 individuals.

276 One of the strengths of our method is the access to a significant history of electronic
277 health records that allow us to derive estimates informed by a greater group of patients
278 throughout the life course. Existing scores^{26,35} imply that the levels of covariates will stay fixed
279 over the life course or require recalculation, which ignores the information within transitions
280 through the life course. Here, our longitudinal outlooks ability allows for individuals to be

281 followed over a lifetime and quickly estimates what their updated risk trajectory would look like
282 under an alternative profile.

283 Estimation of remaining lifetime risk is conducted using age-specific predictions informed
284 only by individuals in the at-risk set at a given age, thus making this a true lifetime estimate. In
285 our work, we choose a conservatively estimated age of 80 as the maximum lifetime age given
286 the density of age estimation with our set. This estimation is possible under the assumption that
287 risk trajectory is similar across shifting windows of age at risk but falls apart with strong calendar
288 time trends. Given that our cohort was required to be between 40 and 70 years old in 2006, we
289 reduced the variation in calendar effects.^{5,36}

290 When combined with genetic information, an emphasis on dynamically updated lifetime
291 risk projections can uncover latent risks in seemingly healthy individuals. Determining an
292 appropriate lifetime risk threshold is a laudable goal.^{2,7} Indeed, current guidelines^{12,36} note that
293 genetic risk scores can identify individuals at birth with a high propensity to develop disease, but
294 few approaches have coupled this information with realized risk stages dynamically. As age
295 increases, short-term risk increases, and the remaining lifetime risk is reduced, meaning that a
296 metric focusing on short-term risk will preferentially focus on disease in older individuals,
297 thwarting the efforts of true prevention. It is not enough to increase the lifetime threshold to
298 account for younger individuals as proposed in European Society of Cardiology guidelines;
299 additional years add additional uncertainty, and thus, having tools capable of dynamically
300 incorporating new information over the life course in combination with more comprehensive time
301 assessments is critical to moving prevention forward. We provide an application for individuals
302 to assess risk in real-time for patients and clinicians (**Supp. Fig. 10**;
303 <https://surbut.shinyapps.io/risk/>)

304 In this study, we use a composite of phenotypic codes to define our risk factor states.
305 One of the challenges of developing a lifetime assessment tool surrounds the availability of
306 continuously updated laboratory data. Using EHR data, an unbiased ascertainment of updated
307 biometric variables at uniform intervals is challenging. We added baseline continuous laboratory
308 data from the age of enrollment to our grid search, and this added little to our model (**Supp. Fig.**
309 **11**).

310 A second limitation surrounds the heterogeneity of phenotyping. We define
311 hyperlipidemia and hypertension according to validated diagnostic codes.³⁷ However, there
312 exists heterogeneity in the severity and duration of these conditions. The potential benefit of

313 adding additional states must be balanced with the uncertainty imposed and the reduction in
314 sample size caused by dispersion across grades of each condition. Our model resolves the loss
315 in underlying latent risk that is often erroneously captured in EHR data when an individual's
316 nominal laboratory value falls secondary to medication use.

317 One of the advantages of heterogenous data collection is a wealth of available
318 phenotyping modalities: the UKBB has access through linkages to routinely available national
319 health systems enhanced by self-report and previous records.³⁸ Although not all individuals
320 included had GP data, we demonstrate that the age and prevalence of conditions is
321 homogenous between individuals in the GP subset and otherwise (**Supp. Figs. 1**) and that
322 analysis on this subset alone results in similar model discrimination.

323 Third, the generalizability of our findings may be impacted by study design and sample
324 specificity. The UK Biobank included healthier and less socioeconomically deprived individuals
325 who were predominantly White Europeans living in the United Kingdom.³⁹ Furthermore, given
326 that the minimum age for genotyping was 40 years old, we began our inference for risk
327 modeling at age 40, provided they were captured in the EHR before then. Although individuals
328 who reached age 40 prior to enrollment were appropriately at risk for the primary CAD outcome
329 given their capture in the longitudinal EHR, they were protected from death until the time of
330 enrollment, which may affect estimates related to the competing risk of death. For time-
331 dependent evaluation of our prediction, we conservatively left-censored at age of enrollment
332 to eliminate years protected from death and found that the improvements in discrimination
333 over FRS30RC remained unchanged. We note consistent performance in external validation
334 in the FOS cohort, where all death and CAD events occurred exclusively after enrollment.
335 Finally, our dynamic logistic regression approach can readily be adapted to any population with
336 minimal computational resources, and we provide code to do so.

337 Leveraging a unique resource of genetic and longitudinal clinical data spanning over 80
338 years in nearly 500,000 participants of the UK Biobank prospective cohort study, we develop
339 MSGene, a multistate model for dynamic transitions throughout the life course to estimate
340 lifetime risk of CAD. MSGene is well-calibrated and discriminates early and late events both in
341 the UK Biobank and an external validation sample. We anticipate that by providing interpretable
342 and dynamic estimates of CAD lifetime risk, MSGene may inform future therapeutic decisions to
343 enable more efficient and effective CAD prevention throughout the life course.

344

345 **Acknowledgments**

346 We would like to acknowledge Leslie Gaffney of the MIT-Broad Communications Lab for her
347 invaluable graphics and copyediting advice.

348

349 **SOURCES OF FUNDING**

350

351 S.M.U. is supported by T32HG010464 from the National Human Genome Research Institute.

352 S.K is supported by the NIH (K23HL169839) and the American Heart Association

353 (23CDA1050571).

354 S.J.C. is supported by a grant of the Korea Health Technology R&D Project through the Korea

355 Health Industry Development Institute (KHIDI), funded by the Ministry of Health and Welfare,

356 Republic of Korea (grant no.: HI19C1330). A.C.F is supported by grants 1K08HL161448 and

357 R01HL164629 from the National Heart, Lung, and Blood Institute.

358 P.T.E reported receiving grants from the NIH (1R01HL092577, 1R01HL157635, and

359 5R01HL139731), the American Heart Association Strategically Focused Research Networks

360 (18SFRN34110082), the European Union (MAESTRIA 965286), Bayer AG (to the Broad

361 Institute), IBM Health (to the Broad Institute), Bristol Myers Squibb (to Massachusetts General

362 Hospital), and Pfizer (to Massachusetts General Hospital).

363 A.G. is supported by National Institutes of Health (NIH) grant nos R01CA227237,

364 R01CA244569 and R21HG010748, and awards from the Claudia Adams Barr Foundation,

365 Louis B. Mayer Foundation, Doris Duke Charitable Foundation, Emerson Collective and Phi

366 Beta Psi Sorority.

367 P.N. is supported by grants R01HL1427, R01HL148565, and R01HL148050 from the National

368 Heart, Lung, and Blood Institute, and grant 1U01HG011719 from the National Human Genome

369 Research Institute.

370

371 **DISCLOSURES**

372

373 During the course of the project, M.W.Y. became a full-time employee of GSK.

374 A.C.F. is co-founder of Goodpath.

375 PTE reports personal fees from Bayer AG, Novartis, and MyoKardia.

376 GP holds equity in Phaeno Biotechnologies, is on the SAB of RealmIDX and currently consults

377 for Delphi Diagnostics.

378 P.N. reports research grants from Allelica, Apple, Amgen, Boston Scientific, Genentech / Roche,

379 and Novartis, personal fees from Allelica, Apple, AstraZeneca, Blackstone Life Sciences,

380 Foresite Labs, Genentech / Roche, GV, HeartFlow, Magnet Biomedicine, and Novartis, scientific

381 advisory board membership of Esperion Therapeutics, Preciseli, and TenSixteen Bio, scientific

382 co-founder of TenSixteen Bio, equity in MyOme, Preciseli, and TenSixteen Bio, and spousal

383 employment at Vertex Pharmaceuticals, all unrelated to the present work. The remaining

384 authors have nothing to disclose.

385

386 Works Cited

- 387 1. Tsao CW, Aday AW, Almarzooq ZI, Anderson CAM, Arora P, Avery CL, Baker-Smith CM,
388 Beaton AZ, Boehme AK, Buxton AE, Commodore-Mensah Y, Elkind MSV, Evenson KR,
389 Eze-Nliam C, Fugar S, Generoso G, Heard DG, Hiremath S, Ho JE, Kalani R, Kazi DS, Ko
390 D, Levine DA, Liu J, Ma J, Magnani JW, Michos ED, Mussolino ME, Navaneethan SD,
391 Parikh NI, Poudel R, Rezk-Hanna M, Roth GA, Shah NS, St-Onge M-P, Thacker EL, Virani
392 SS, Voeks JH, Wang N-Y, Wong ND, Wong SS, Yaffe K, Martin SS, Subcommittee on
393 behalf of the AHAC on E and PSC and SS. Heart Disease and Stroke Statistics—2023
394 Update: A Report From the American Heart Association. *Circulation* [Internet]. 2023 [cited
395 2023 May 20]; Available from:
396 <https://www.ahajournals.org/doi/abs/10.1161/CIR.0000000000001123>
- 397 2. Lloyd-Jones DM, Leip EP, Larson MG, D'Agostino RB, Beiser A, Wilson PWF, Wolf PA,
398 Levy D. Prediction of Lifetime Risk for Cardiovascular Disease by Risk Factor Burden at 50
399 Years of Age. *Circulation*. 2006;113:791–798.
- 400 3. Wilkins JT, Karmali KN, Huffman MD, Allen NB, Ning H, Berry JD, Garside DB, Dyer A,
401 Lloyd-Jones DM. Data Resource Profile: The Cardiovascular Disease Lifetime Risk
402 Pooling Project. *Int J Epidemiol*. 2015;44:1557–1564.
- 403 4. Bundy JD, Ning H, Zhong VW, Paluch AE, Lloyd-Jones DM, Wilkins JT, Allen NB.
404 Cardiovascular Health Score and Lifetime Risk of Cardiovascular Disease. *Circulation:*
405 *Cardiovascular Quality and Outcomes* [Internet]. 2020 [cited 2023 Jun 13]; Available from:
406 <https://www.ahajournals.org/doi/abs/10.1161/CIRCOUTCOMES.119.006450>
- 407 5. Grundy SM, Stone NJ, Bailey AL, Beam C, Birtcher KK, Blumenthal RS, Braun LT, de
408 Ferranti S, Faiella-Tommasino J, Forman DE, Goldberg R, Heidenreich PA, Hlatky MA,
409 Jones DW, Lloyd-Jones D, Lopez-Pajares N, Ndumele CE, Orringer CE, Peralta CA,
410 Saseen JJ, Smith SC, Sperling L, Virani SS, Yeboah J. 2018
411 AHA/ACC/AACVPR/AAPA/ABC/ACPM/ADA/AGS/ APhA/ASPC/NLA/PCNA Guideline on
412 the Management of Blood Cholesterol: Executive Summary. *Circulation*. 2019;139:e1082–
413 e1143.
- 414 6. Yadlowsky S, Hayward RA, Sussman JB, McClelland RL, Min Y-I, Basu S. Clinical
415 Implications of Revised Pooled Cohort Equations for Estimating Atherosclerotic
416 Cardiovascular Disease Risk. *Ann Intern Med*. 2018;169:20–29.
- 417 7. Navar AM, Fine LJ, Ambrosius WT, Brown A, Douglas PS, Johnson K, Khera AV, Lloyd-
418 Jones D, Michos ED, Mujahid M, Muñoz D, Nasir K, Redmond N, Ridker PM, Robinson J,
419 Schopfer D, Tate DF, Lewis CE. Earlier treatment in adults with high lifetime risk of
420 cardiovascular diseases: What prevention trials are feasible and could change clinical
421 practice? Report of a National Heart, Lung, and Blood Institute (NHLBI) Workshop.
422 *American Journal of Preventive Cardiology*. 2022;12:100430.
- 423 8. Jaspers NEM, Blaha MJ, Matsushita K, van der Schouw YT, Wareham NJ, Khaw K-T,
424 Geisel MH, Lehmann N, Erbel R, Jöckel K-H, van der Graaf Y, Verschuren WMM, Boer
425 JMA, Nambi V, Visseren FLJ, Dorresteyn JAN. Prediction of individualized lifetime benefit

- 426 from cholesterol lowering, blood pressure lowering, antithrombotic therapy, and smoking
427 cessation in apparently healthy people. *Eur Heart J*. 2020;41:1190–1199.
- 428 9. Navar AM, Fonarow GC, Pencina MJ. Time to Revisit Using 10-Year Risk to Guide Statin
429 Therapy. *JAMA Cardiol*. 2022;7:785.
- 430 10. Zeitouni M, Nanna MG, Sun J-L, Chiswell K, Peterson ED, Navar AM. Performance of
431 Guideline Recommendations for Prevention of Myocardial Infarction in Young Adults.
432 *Journal of the American College of Cardiology*. 2020;76:653–664.
- 433 11. Lloyd-Jones DM, Albert MA, Elkind M. The American Heart Association’s Focus on
434 Primordial Prevention. *Circulation*. 2021;144:e233–e235.
- 435 12. 2021 ESC Guidelines on cardiovascular disease prevention in clinical practice | European
436 Heart Journal | Oxford Academic [Internet]. [cited 2023 Oct 6]; Available from:
437 <https://academic.oup.com/eurheartj/article/42/34/3227/6358713>
- 438 13. Berry JD, Dyer A, Cai X, Garside DB, Ning H, Thomas A, Greenland P, Van Horn L, Tracy
439 RP, Lloyd-Jones DM. Lifetime Risks of Cardiovascular Disease. *New England Journal of*
440 *Medicine*. 2012;366:321–329.
- 441 14. Conner SC, Beiser A, Benjamin EJ, LaValley MP, Larson MG, Trinquart L. A comparison
442 of statistical methods to predict the residual lifetime risk. *Eur J Epidemiol*. 2022;37:173–
443 194.
- 444 15. Michos ED, Choi AD. Coronary Artery Disease in Young Adults. *Journal of the American*
445 *College of Cardiology*. 2019;74:1879–1882.
- 446 16. O’Sullivan JW, Raghavan S, Marquez-Luna C, Luzum JA, Damrauer SM, Ashley EA,
447 O’Donnell CJ, Willer CJ, Natarajan P, on behalf of the American Heart Association Council
448 on Genomic and Precision Medicine; Council on Clinical Cardiology; Council on
449 Arteriosclerosis, Thrombosis and Vascular Biology; Council on Cardiovascular Radiology
450 and Intervention; Council on Lifestyle and Cardiometabolic Health; and Council on
451 Peripheral Vascular Disease. Polygenic Risk Scores for Cardiovascular Disease: A
452 Scientific Statement From the American Heart Association. *Circulation* [Internet]. 2022
453 [cited 2023 Oct 6];146. Available from:
454 <https://www.ahajournals.org/doi/10.1161/CIR.0000000000001077>
- 455 17. Inouye M, Abraham G, Nelson CP, Wood AM, Sweeting MJ, Dudbridge F, Lai FY, Kaptoge
456 S, Brozynska M, Wang T, Ye S, Webb TR, Rutter MK, Tzoulaki I, Patel RS, Loos RJF,
457 Keavney B, Hemingway H, Thompson J, Watkins H, Deloukas P, Di Angelantonio E,
458 Butterworth AS, Danesh J, Samani NJ. Genomic Risk Prediction of Coronary Artery
459 Disease in 480,000 Adults. *Journal of the American College of Cardiology*. 2018;72:1883–
460 1893.
- 461 18. Sniderman AD, Furberg CD. Age as a modifiable risk factor for cardiovascular disease.
462 *Lancet*. 2008;371:1547–1549.

- 463 19. Wang N, Woodward M, Huffman MD, Rodgers A. Compounding Benefits of Cholesterol-
464 Lowering Therapy for the Reduction of Major Cardiovascular Events: Systematic Review
465 and Meta-Analysis. *Circulation: Cardiovascular Quality and Outcomes*. 2022;15:e008552.
- 466 20. Marma AK, Berry JD, Ning H, Persell SD, Lloyd-Jones DM. Distribution of 10-year and
467 lifetime predicted risks for cardiovascular disease in US adults: findings from the National
468 Health and Nutrition Examination Survey 2003 to 2006. *Circ Cardiovasc Qual Outcomes*.
469 2010;3:8–14.
- 470 21. Le-Rademacher JG, Therneau TM, Ou F-S. The Utility of Multistate Models: A Flexible
471 Framework for Time-to-Event Data. *Curr Epidemiol Rep*. 2022;9:183–189.
- 472 22. Wreede LC de, Fiocco M, Putter H. **mstate**: An R Package for the Analysis of Competing
473 Risks and Multi-State Models. *J Stat Soft* [Internet]. 2011 [cited 2022 Dec 1];38. Available
474 from: <http://www.jstatsoft.org/v38/i07/>
- 475 23. Brookmeyer R, Abdalla N. Multistate models and lifetime risk estimation: Application to
476 Alzheimer’s disease. *Statistics in Medicine*. 2019;38:1558–1565.
- 477 24. Neumann JT, Thao LTP, Callander E, Carr PR, Qaderi V, Nelson MR, Reid CM, Woods
478 RL, Orchard SG, Wolfe R, Polekhina G, Williamson JD, Trauer JM, Newman AB, Murray
479 AM, Ernst ME, Tonkin AM, McNeil JJ. A multistate model of health transitions in older
480 people: a secondary analysis of ASPREE clinical trial data. *The Lancet Healthy Longevity*.
481 2022;3:e89–e97.
- 482 25. Jack CR, Therneau TM, Wiste HJ, Weigand SD, Knopman DS, Lowe VJ, Mielke MM,
483 Vemuri P, Roberts RO, Machulda MM, Senjem ML, Gunter JL, Rocca WA, Petersen RC.
484 Rates of transition between amyloid and neurodegeneration biomarker states and to
485 dementia among non-demented individuals: a population-based cohort study. *Lancet*
486 *Neurol*. 2016;15:56–64.
- 487 26. Pencina MJ, Ralph B, D’Agostino S, Larson MG, Massaro JM, Vasan RS. Predicting the
488 30-Year Risk of Cardiovascular Disease. *Circulation* [Internet]. 2009 [cited 2023 Sep
489 20]; Available from:
490 <https://www.ahajournals.org/doi/abs/10.1161/CIRCULATIONAHA.108.816694>
- 491 27. Cleveland WS. Robust Locally Weighted Regression and Smoothing Scatterplots. *Journal*
492 *of the American Statistical Association*. 1979;74:829–836.
- 493 28. Rospleszcz S, Starnecker F, Linkohr B, von Scheidt M, Gieger C, Schunkert H, Peters A.
494 Validation of the 30-Year Framingham Risk Score in a German Population-Based Cohort.
495 *Diagnostics (Basel)*. 2022;12:965.
- 496 29. Cholesterol Treatment Trialists’ (CTT) Collaborators, Mihaylova B, Emberson J, Blackwell
497 L, Keech A, Simes J, Barnes EH, Voysey M, Gray A, Collins R, Baigent C. The effects of
498 lowering LDL cholesterol with statin therapy in people at low risk of vascular disease:
499 meta-analysis of individual data from 27 randomised trials. *Lancet*. 2012;380:581–590.
- 500 30. Cholesterol Treatment Trialists’ (CTT) Collaborators, Mihaylova B, Emberson J, Blackwell
501 L, Keech A, Simes J, Barnes EH, Voysey M, Gray A, Collins R, Baigent C. The effects of

- 502 lowering LDL cholesterol with statin therapy in people at low risk of vascular disease:
503 meta-analysis of individual data from 27 randomised trials. *Lancet*. 2012;380:581–590.
- 504 31. Chou R, Cantor A, Dana T, Wagner J, Ahmed AY, Fu R, Ferencik M. Statin Use for the
505 Primary Prevention of Cardiovascular Disease in Adults: Updated Evidence Report and
506 Systematic Review for the US Preventive Services Task Force. *JAMA*. 2022;328:754.
- 507 32. Natarajan P, Young R, Stitzel NO, Padmanabhan S, Baber U, Mehran R, Sartori S, Fuster
508 V, Reilly DF, Butterworth A, Rader DJ, Ford I, Sattar N, Kathiresan S. Polygenic Risk
509 Score Identifies Subgroup With Higher Burden of Atherosclerosis and Greater Relative
510 Benefit From Statin Therapy in the Primary Prevention Setting. *Circulation*.
511 2017;135:2091–2101.
- 512 33. Thanassoulis G, Sniderman AD, Pencina MJ. A Long-term Benefit Approach vs Standard
513 Risk-Based Approaches for Statin Eligibility in Primary Prevention. *JAMA Cardiol*.
514 2018;3:1090–1095.
- 515 34. Pencina MJ, Pencina KM, Lloyd-Jones D, Catapano AL, Thanassoulis G, Sniderman AD.
516 The Expected 30-Year Benefits of Early Versus Delayed Primary Prevention of
517 Cardiovascular Disease by Lipid Lowering. *Circulation*. 2020;142:827–837.
- 518 35. Hippisley-Cox J, Coupland C, Vinogradova Y, Robson J, Minhas R, Sheikh A, Brindle P.
519 Predicting cardiovascular risk in England and Wales: prospective derivation and validation
520 of QRISK2. *BMJ*. 2008;336:1475–1482.
- 521 36. Arnett DK, Blumenthal RS, Albert MA, Buroker AB, Goldberger ZD, Hahn EJ, Himmelfarb
522 CD, Khera A, Lloyd -Jones Donald, McEvoy JW, Michos ED, Miedema MD, Mu ñoz D,
523 Smith SC, Virani SS, Williams KA, Yeboah J, Ziaieian B. 2019 ACC/AHA Guideline on the
524 Primary Prevention of Cardiovascular Disease. *Journal of the American College of*
525 *Cardiology*. 2019;74:e177–e232.
- 526 37. Yeung MW, Van Der Harst P, Verweij N. ukbpheno v1.0: An R package for phenotyping
527 health-related outcomes in the UK Biobank. *STAR Protocols*. 2022;3:101471.
- 528 38. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, Downey P, Elliott P, Green J,
529 Landray M, Liu B, Matthews P, Ong G, Pell J, Silman A, Young A, Sprosen T, Peakman T,
530 Collins R. UK Biobank: An Open Access Resource for Identifying the Causes of a Wide
531 Range of Complex Diseases of Middle and Old Age. *PLoS Med*. 2015;12:e1001779.
- 532 39. Fry A, Littlejohns TJ, Sudlow C, Doherty N, Adamska L, Sprosen T, Collins R, Allen NE.
533 Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank
534 Participants With Those of the General Population. *American Journal of Epidemiology*.
535 2017;186:1026–1034.
- 536 40. Klarin D, Zhu QM, Emdin CA, Chaffin M, Horner S, McMillan BJ, Leed A, Weale ME,
537 Spencer CCA, Aguet F, Segrè AV, Ardlie KG, Khera AV, Kaushik VK, Natarajan P,
538 CARDIoGRAMplusC4D Consortium, Kathiresan S. Genetic analysis in UK Biobank links
539 insulin resistance and transendothelial migration pathways to coronary artery disease. *Nat*
540 *Genet*. 2017;49:1392–1397.

- 541 41. Thompson DJ, Wells D, Selzam S, Peneva I, Moore R, Sharp K, Tarran WA, Beard EJ,
542 Riveros-Mckay F, Palmer D, Seth P, Harrison J, Futema M, Consortium GER, McVean G,
543 Plagnol V, Donnelly P, Weale ME. UK Biobank release and systematic evaluation of
544 optimised polygenic risk scores for 53 diseases and quantitative traits [Internet]. 2022
545 [cited 2023 Oct 2];2022.06.16.22276246. Available from:
546 <https://www.medrxiv.org/content/10.1101/2022.06.16.22276246v1>
- 547 42. Darke P, Cassidy S, Catt M, Taylor R, Missier P, Bacardit J. Curating a longitudinal
548 research resource using linked primary care EHR data—a UK Biobank case study. *Journal*
549 *of the American Medical Informatics Association*. 2022;29:546–552.
- 550 43. Cleveland WS, Devlin SJ. Locally Weighted Regression: An Approach to Regression
551 Analysis by Local Fitting.
- 552 44. Ference BA, Yoo W, Alesh I, Mahajan N, Mirowska KK, Mewada A, Kahn J, Afonso L,
553 Williams KA, Flack JM. Effect of long-term exposure to lower low-density lipoprotein
554 cholesterol beginning early in life on the risk of coronary heart disease: a Mendelian
555 randomization analysis. *J Am Coll Cardiol*. 2012;60:2631–2639.
- 556 45. Ference BA. How to use Mendelian randomization to anticipate the results of randomized
557 trials. *European Heart Journal*. 2018;39:360–362.
- 558 46. Mega J, Stitziel N, Smith J, Chasman D, Caulfield M, Devlin J, Nordio F, Hyde C, Cannon
559 C, Sacks F, Poulter N, Sever P, Ridker P, Braunwald E, Melander O, Kathiresan S,
560 Sabatine M. Genetic Risk, Coronary Heart Disease Events, and the Clinical Benefit of
561 Statin Therapy. *Lancet*. 2015;385:2264–2271.
- 562 47. Marston NA, Kamanu FK, Nordio F, Gurmu Y, Roselli C, Sever PS, Pedersen TR, Keech
563 AC, Wang H, Lira Pineda A, Giugliano RP, Lubitz SA, Ellinor PT, Sabatine MS, Ruff CT.
564 Predicting Benefit From Evolocumab Therapy in Patients With Atherosclerotic Disease
565 Using a Genetic Risk Score: Results From the FOURIER Trial. *Circulation*. 2020;141:616–
566 623.
- 567 48. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D,
568 Delaneau O, O’Connell J, Cortes A, Welsh S, Young A, Effingham M, McVean G, Leslie S,
569 Allen N, Donnelly P, Marchini J. The UK Biobank resource with deep phenotyping and
570 genomic data. *Nature*. 2018;562:203–209.
- 571 49. Therneau T, Crowson C, Atkinson E. Using Time Dependent Covariates and Time
572 Dependent Coefficients in the Cox Model. :31.
- 573 50. Harrell FE Jr, Califf RM, Pryor DB, Lee KL, Rosati RA. Evaluating the Yield of Medical
574 Tests. *JAMA*. 1982;247:2543–2546.
- 575 51. Feinleib M, Kannel WB, Garrison RJ, McNamara PM, Castelli WP. The Framingham
576 Offspring Study. Design and preliminary data. *Prev Med*. 1975;4:518–525.
- 577
- 578

579
580

Figure Legends

581 **Figure 1. Multistate transitions over time.**

582 **A.** We depict the potential one-step transitions in our multistate framework. Per year, an
583 individual can progress from health to single risk factor states, CAD or death. Similarly, an
584 individual can progress from single risk factor states, to double risk factor states, to CAD or
585 death; from double risk factor states, to triple risk factor, CAD or death. **B.** We display the
586 proportional occupancy excluding censored individuals at each state.

587 **CAD:** coronary artery disease, **Ht:** hypertension, **HyperLip:** hyperlipidemia, **Dm:** Type 2
588 diabetes mellitus.

589

590 **Figure 2. Study overview.**

591 **A.** Using the UK Biobank data on half a million participants (54% female) with access to health
592 record from 1940, we harmonize hospitalization, prescription and primary care records from the
593 EHR and train our model on individuals free of CAD at age 40. The UKB required participants to
594 be between ages 40–69 between 2006–2010 for genotyping. In our model, individuals join
595 disease-free in the ‘health’ state and progress to additional states upon censoring. We use 80%
596 of the eligible data for training and the remaining 20% for testing. For the testing subset we
597 require that individuals have variables necessary for computation of FRS30 (and FRS30RC)
598 and the pooled cohort equations, which require laboratory (HDL, TC) and biometric (SBP)
599 measurements. **B.** For a sample patient, we document the construction of our cohort. This
600 individual is first observed in the health record at age 25; he is diagnosed with hypertension at
601 age 39, and begins informing our risk estimation for CAD at age 40 in the hypertensive
602 category. He transitions to the hypertension and hyperlipidemia category at age 50, 25 years
603 after first encounter and 10 years after entering our risk estimation, thus contributing 10 years of
604 data.

605 **TC:** total cholesterol, **SBP:** systolic blood pressure, **HDL:** high-density lipoprotein, **CAD:**
606 coronary artery disease, **FRS30:** Framingham 30 year, **FRS30RC:** Framingham 30 year
607 recalibrated, **PCE:** Pooled cohort equation 10-year risk; **EHR:** electronic health record.

608

609

610 **Figure 3. Survival, 10-year and lifetime risk curves.**

611 In **A.**, we demonstrate the singular projected survival curve by MSGene for an individual at age
612 40 of low, medium or high genomic risk. In **B.** we demonstrate the MSGene predicted 10-year

613 risk for individuals at each age along the x-axis, showing that, in general, for fixed window
614 approaches, 10-year risk is monotonically increasing. In **C**, we demonstrate the MSGene
615 predicted lifetime risk curve for individuals at each age featured along the x-axis under an
616 untreated (dashed) or treated (solid) strategy. The conditional remaining lifetime risk declines
617 with age, from 24% for a high genomic risk individual in our cohort to <5% for an individual at
618 the same risk level by age 70. In **D**, using the FRS30RC equation, like 10-year risk and unlike
619 the remaining lifetime risk approach, 30-year risk calculation is monotonically increasing, from
620 13.4 (13.2–13.6%) at age 40 to 32.9% at age 70 for an individual of the highest genomic risk.
621 **FRS30RC**: Framingham 30 year recalibrated.

622

623 **Figure 4: Time-dependent threshold analysis.**

624 We consider the distribution of the first age at which an individual exceeds the PCE-derived 10-
625 year threshold of 5% (**A**), or lifetime threshold or 10% using FRS30RC (**B**) or the MSGene
626 lifetime prediction (**C**). We then use this age as a time-dependent predictor of time-to-event in a
627 time-dependent Cox PH (**Supp. methods**) in which an individual's time followed is stratified by
628 start time and periods in which a threshold is passed, and final censoring time with an indicator
629 variable demarcating whether or not each threshold has been surpassed. We left censor these
630 intervals at age of enrollment conservatively to exclude time protected from death. We report
631 Harrell's C-index ($p < 2 \times 10^{-16}$) for discrimination on how well a model predicts events that tend
632 to occur earlier versus later. Left-facing indicate individuals who surpass the threshold at first
633 prediction, and right-facing arrow indicates individuals who never surpass a threshold for a
634 given metric. FRS30RC is shown here with C-index 0.52 (original FRS30 C-index 0.50) vs.
635 MSGene 0.72, $p < 2 \times 10^{-16}$ (**D**). We compute the lifetime prediction at each age under one of
636 eight potential risk starting states, with bootstrapped confidence intervals for a sample individual
637 (**E**). Using the electronic health record, we extract state position for each individual per year. We
638 then use MSGene to compute predicted risk for each individual at each state in time, displayed
639 here for a sample individuals (**F**). We use these as predictors in a time-dependent Cox model in
640 which we expand the data set into non-overlapping intervals for each individual (**Supp.**
641 **methods; Supp. Fig. 17**) and conservatively left censor before enrollment to avoid time
642 protected from death. We evaluate the concordance when compared to FRS30RC and PCE-
643 derived 10-year, $p < 2.2 \times 10^{-16}$ (**G**).

644 **FRS30RC**: Framingham 30-year recalibrated, **PCE**: pooled cohort equations, **Cox PH**: Cox
645 proportional hazards model

646

647 **Figure 5: Absolute risk reduction: Short-term and lifetime risk.**

648 We display the relationship between remaining lifetime and 10-year risk. Each ray represents an
649 age group, in which individuals are parameterized by their short- (10-year) and long-term
650 (lifetime) risk, and colored by genomic risk in SD from mean. We display the lifetime absolute
651 risk reduction as computed in Equation RR and stratified by age rays, and colored by genetic
652 risk. **(A)** For an individual at the top genetic risk at age 40, MSGene predicted 10-year risk is
653 roughly equivalent to an individual at the lowest genetic risk at age 70 (3.8% vs 4.2%, SE 0.01).
654 However, the MSGene projected lifetime benefit is directly proportional to lifetime risk **(B)**, and
655 more than twice that of a high risk individual at age 70 (5.0 vs 2.3%, SEM 0.02). **(C)**
656 Marginalized across starting states and covariate profiles, we project absolute risk difference
657 (%) under a treated and untreated setting. At age 40, this ranges from a median of 5.8% (SD
658 0.01) to 0.8% (SD 0.01) at age 79.

659 **SEM:** standard error of mean, **RR:** relative risk, **SD:** CAD-PRS SD.

	Low (N=96235)	Intermediate (N=288563)	High (N=95840)	Overall (N=480638)
Sex				
Female	51958 (54.0%)	156570 (54.3%)	52125 (54.4%)	260653 (54.2%)
Male	44277 (46.0%)	131993 (45.7%)	43715 (45.6%)	219985 (45.8%)
Birthdate				
Median [Min, Max]	1950 [1940, 1970]	1950 [1930, 1970]	1950 [1940, 1970]	1950 [1930, 1970]
Age First Enrolled in NHS				
Mean (SD)	29.2 (13.2)	29.2 (13.2)	29.1 (13.2)	29.2 (13.2)
Median [Min, Max]	24.5 [18.0, 78.6]	24.3 [18.0, 78.3]	24.2 [18.0, 79.1]	24.3 [18.0, 79.1]
Years Followed				
Mean (SD)	29.5 (8.05)	29.5 (8.03)	29.3 (8.01)	29.4 (8.03)
Median [Min, Max]	30.6 [0.375, 45.5]	30.6 [0.843, 47.6]	30.3 [1.36, 44.8]	30.5 [0.375, 47.6]
Develop Hypertension				
No	63687 (66.2%)	174497 (60.5%)	52002 (54.3%)	290186 (60.4%)
Yes	32548 (33.8%)	114066 (39.5%)	43838 (45.7%)	190452 (39.6%)
Develop Coronary Disease				
No	89929 (93.4%)	258215 (89.5%)	79034 (82.5%)	427178 (88.9%)
Yes	6306 (6.6%)	30348 (10.5%)	16806 (17.5%)	53460 (11.1%)
Develop Hyperlipidemia				
No	79046 (82.1%)	221300 (76.7%)	66698 (69.6%)	367044 (76.4%)
Yes	17189 (17.9%)	67263 (23.3%)	29142 (30.4%)	113594 (23.6%)
Current Smoker				
No	86517 (89.9%)	258134 (89.5%)	85315 (89.0%)	429966 (89.5%)
Yes	9718 (10.1%)	30429 (10.5%)	10525 (11.0%)	50672 (10.5%)
Proportion White				
Yes	82842 (86.1%)	251780 (87.3%)	82479 (86.1%)	417101 (86.8%)
General Practice Registry Members				
Not Member	52539 (54.6%)	155429 (53.9%)	51319 (53.5%)	259287 (53.9%)
Member	43696 (45.4%)	133134 (46.1%)	44521 (46.5%)	221351 (46.1%)

660

661 **Table 1. Distribution of overall cohort.** We use approximately 80% (385,541) individuals in
662 the training, and 79,119 in the testing set, of which approximately 45% represent members of
663 the general practice primary care data. Of note, low genomic risk connotes individuals in the
664 lowest (<20%) of genomic risk by PRS percentile, intermediate (20–80%) PRS percentile, and
665 high denotes >80% PRS percentile.

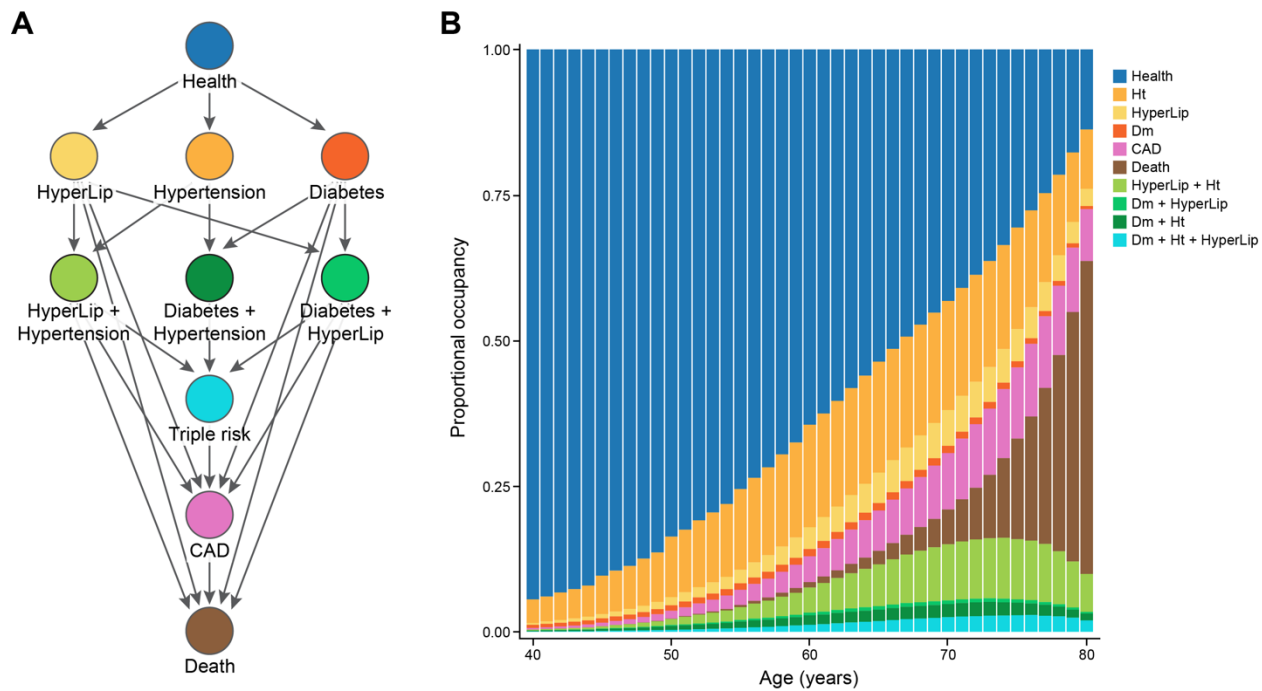
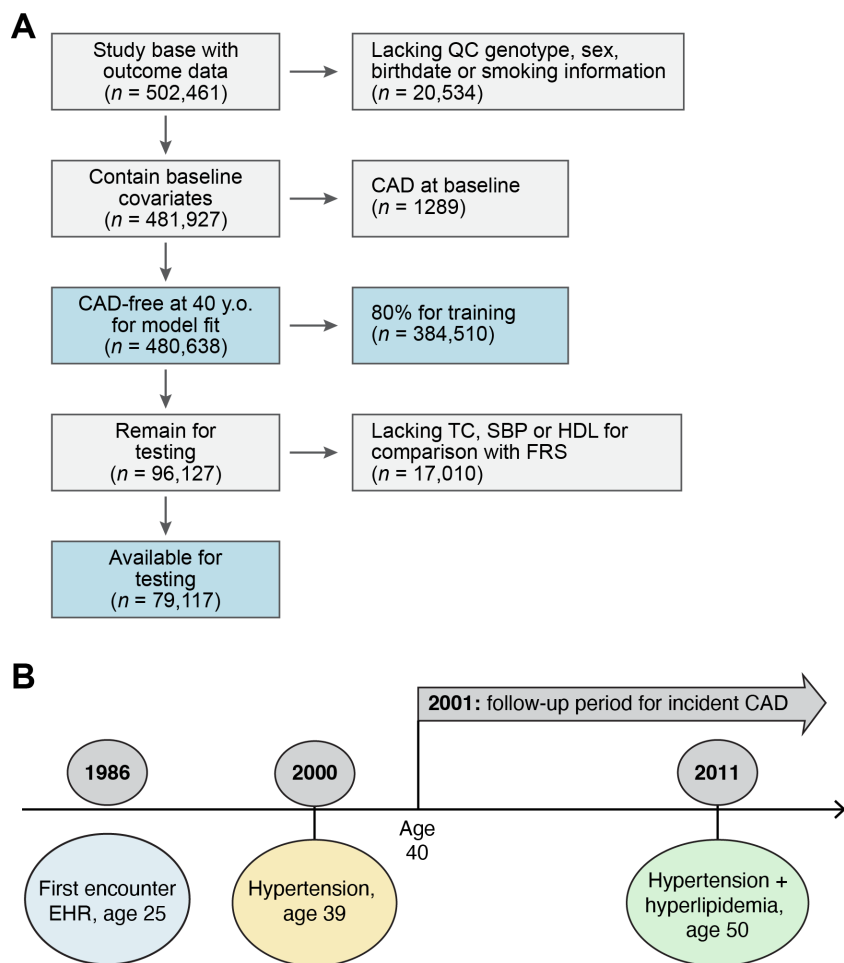


Figure 1.

668



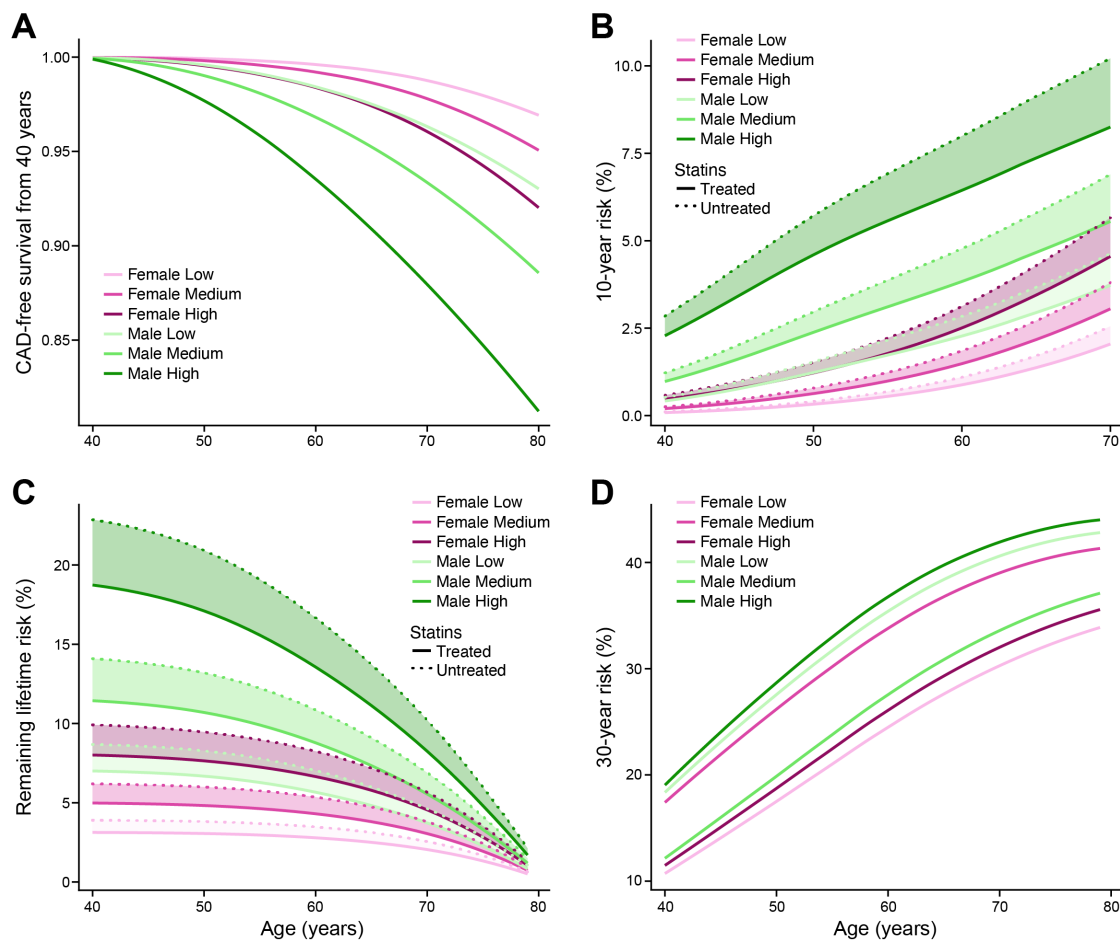
669

670

671

Figure 2.

672



673

674

675

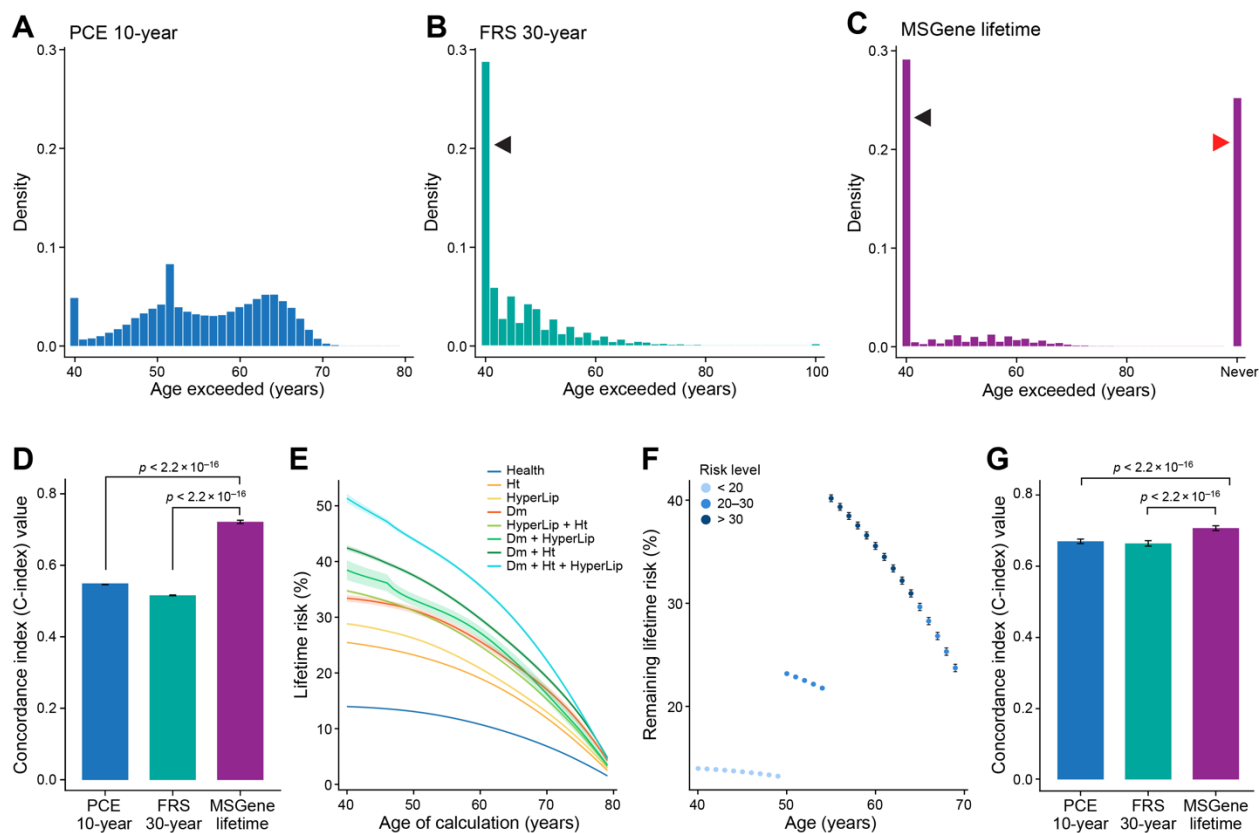
676

677

678

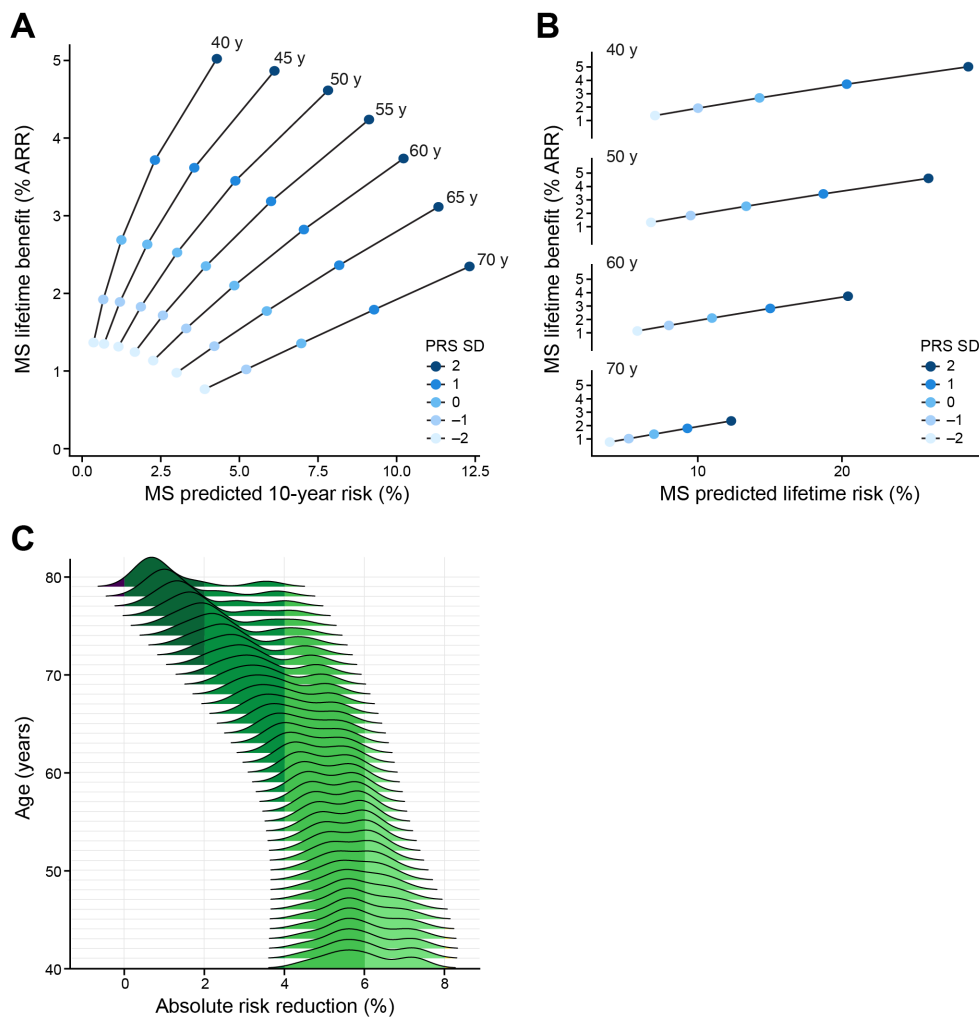
679

Figure 3.



680
681
682

Figure 4.



683
684
685
686

Figure 5.

687 **Methods**

688 *Data source*

689 The UK Biobank (UKB) is a prospective UK population-based study that enrolled approximately
690 half a million adults aged 40–70 between 2006 and 2010 designed to investigate the genetic
691 and lifestyle determinants for a wide range of diseases. Participants underwent genome-wide
692 genotyping, with linkage to longitudinal hospitalization, primary care (GP), and self-report data
693 dating back to 1940 (**Fig. 2; Supp. Figs. 12-13**).³⁷ Using the *ukbpheno* package (version 1.0),³⁷
694 we assembled detailed longitudinal data from the various sources documenting events from
695 1940 until December 2021 for 481,927 individuals after excluding 20,534 who lacked quality
696 control genotyping or risk factor information (**Fig. 2; Supp. Fig. 12-14**). At the time of analysis,
697 linkage to the United Kingdom General Practice (GP) Registry was available for a subset of
698 221,351 individuals. This assembly across data-sources generated phenotypes for hypertension
699 (Htn), diabetes mellitus (DM) (type 1 or 2), hyperlipidemia (Hld), or coronary artery disease
700 (CAD) based on validated collections of hospitalization (HESIN), diagnostic, operation, general
701 practice (GP) clinical and script as well as death information.³⁷ We found high overlap between
702 these phenotypes and our own lab’s previously generated HESIN-restricted phenotypes^{32,40}
703 (**Supp. Fig. 14**). These phenotypes subsequently became the risk factor states in our model.
704 Informed consent was obtained from all participants, and secondary data analyses were
705 approved by the Mass General Brigham Institutional Review Board 2021P002228. Secondary
706 data analysis of UKB was performed under application number 7089.

707 Because of the longitudinal nature of this cohort, every individual is observed at first
708 encounter with the electronic health record (EHR) in early adulthood (median age 24.2 years).
709 We selected UKB participants free of CAD at age 40 and followed until the occurrence of CAD,
710 death, or loss to follow-up (median follow-up 29.9 years). We categorize individuals by their
711 condition at entry into our cohort at age 40 provided they have been observed in the EHR (**Fig.**
712 **2**). We then re-evaluate at each age the risk set as those individuals who have 1) been
713 observed and 2) have not been censored for a given phenotype. We demonstrate the diversity
714 of data sources and the corresponding availability of each data source over time for all
715 considered phenotypes (**Supp. Fig. 13**). In general, our model allows for the progression from
716 CAD to death, but we report here the risk of progression to CAD on CAD-free individuals at
717 baseline.

718 **Polygenic risk**

719 An additional novelty of our model is the incorporation of the dynamic effects of genetics over
720 time. We use CAD polygenic risk score (PRS) as released through the UKB resource⁴¹ and
721 compute on individuals with adequate genotype information after quality control and after
722 controlling for the principal component axes obtained from the common genotype data in the
723 1000 Genomes reference data set using standard methods⁴¹. Data supporting these scores
724 were entirely from external GWAS data (the Standard PRS set) as conducted by Genomics PLC
725 (Oxford, UK) under UKB project 9659.⁴¹ We demonstrate that the distribution of PRS is similar
726 across entry age (**Supplementary Figure 15**).

727 **Statistical analysis**

728 **Detailed Equations**

729 Let π_{jkia} represent the annual transition probability from state **j** to state **k** for individual **i** during
730 year **a**. We let the states **j** and **k** represent time-dependent phenotypes ascertained from the
731 electronic health record such that every individual is in the at-risk 'healthy' set until first
732 censoring. For **p**-covariates for a given individual transitioning from state **j** to **k** we refer to the
733 following equation. 'From' states **J** include Health; single risk factor states: Hypertension (Ht),
734 Hyperlipidemia (Hld), Diabetes Mellitus Type 1 and Type 2 (DM), double risk factor states: Ht &
735 Hld, Ht & Dm, Dm & Hld; Triple risk factor states: Dm & Hld & Ht; and Coronary Artery Disease
736 (CAD). States **K** include all of the 'From' states and Death. For our purposes, we report the
737 progression to CAD or death from any of the starting states included in **J**.

738

$$739 \log \frac{\pi_{jkia}}{1 - \pi_{jkia}} = \hat{\beta}_{jka0} + \hat{\beta}_{jka1}x_1 + \dots + \hat{\beta}_{jkap}x_p$$

740 **Equation 1.**

741

742 Taking the inverse logit of the estimate returns the absolute risk for any individual **i** is a function
743 of the age-specific coefficients and his **p** covariates, such that the annual risk estimate from
744 state **j** to state **k** satisfies:

745

$$746 \pi_{jkia} = \frac{\exp^{X_{ia}B_{jka}}}{1 + \exp^{X_{ia}B_{jka}}}$$

747 **Equation 2.**

748

749 Here we let X represent the $N \times P$ matrix of individuals and covariate profiles at a given age and
750 β represents the $P \times 1$ vector of age and state-state specific smoothed coefficients.

751
752 In equation 2, state j represents the ‘from’ state and state k represents the ‘to’ state. To account
753 for censoring, an individual exits the ‘at risk’ group for transition inference when they are lost to
754 follow-up. We use a one-year interval over which to discretize age intervals and independently
755 estimate the π_{jkia} age-dependent-state to state transitions. We use a limited number of time-
756 fixed covariates: that is sex and polygenic risk score (PRS) and estimate time-dependent
757 effects. We assess current smoker at enrollment in the UK Biobank and use as a time fixed
758 effect for model estimation – that is, individuals reporting ‘current smoker’ at enrollment in the
759 UKBB are considered as smokers in each age-specific logistic regression. For inference of time-
760 dependent covariates, we treat both anti-hypertensive and statin use as individual time-
761 dependent covariate which is reevaluated at each year of model estimation using prescription
762 data from the UKB.⁴² Our final prediction model allows for continuous updates of smoking and
763 medication usage in estimating age-specific transition probabilities. We use 80% of our data as
764 training and 20% as testing (**Fig. 2**) for internal cross-validation and to optimize model fit.
765 Accordingly, this divides our data into a training set for model fitting using 384,510 samples and
766 a testing data set of 79,117 unique individuals.

767 **Predicted Interval Risk**

768 Predicted risk over a given time interval for a given individual i of progressing to state k from
769 state j over any Y -year period ranging from age A_1 to A_2 is where a indexes the current age:
770

$$771 \quad \text{Interval Risk} = 1 - \prod_{A_1}^{A_2} (1 - \pi_{jkia})$$

772 **Equation 3.**

773 Accordingly, risk for an individual i of progressing to state k from state j where L is the
774 maximum age of life and a is the currently observed age. For our purposes, we choose $L = 80$
775 in line with the available data by age in the UK Biobank.

776

$$777 \quad \text{Remaining Lifetime Risk} = 1 - \prod_{A_1}^L (1 - \pi_{jkia})$$

778 **Equation 4.**

779
780
781
782
783
784
785

The remaining lifetime risk can be modified to account for treatments by applying a constant relative risk reduction to the age-specific transition probabilities in expression 4. Then the interval risk under treatment can be calculated using the per-year risk reduction **RR** of progressing to state *k* from state *j* over an interval from age **A₁** to **A₂** is:

786

$$\text{Interval Risk under treatment} = 1 - \prod_{A_1}^{A_2} (1 - (1 - RR) \times \pi_{jkia})$$

787 **Equation 5.**

788

789 For the purposes of this manuscript, we are interested in state *k* = CAD. We impute the relative
790 risk reduction of 0.20 from 24 trials of statin therapy.²⁹ Within our model, we constrain each
791 individual's predicted probabilities across states per year to sum to one such that for each age
792 **a**, the probability of staying within the given state is the complement of the sum of transitions
793 over *K* to the alternative states:

794

$$\pi_{jjia} = 1 - \sum_{k \neq j} \pi_{jkia}$$

795 **Equation 6.**

796 It is somewhat arbitrary to choose *j* as the "to" state whose probability is determined as the
797 complement of the others. We choose *j* because it is mostly above 50% and the constraint in 6
798 will guarantee that for a given age the probabilities for an individual of a particular covariate
799 profile sum to 1. The alternative of fitting a polytomous regression is computationally much more
800 demanding and gives approximately the same answer.

801

802 **Flexible Smoothing Across Ages**

803 We extract the unsmoothed coefficients $\hat{\beta}_{jka}$ for each age and state transition from the logistic
804 regressions in (2). To borrow information across ages, we fit a smoothed locally estimated
805 polynomial regression in which for each state to state transition and each covariate, we fit a
806 locally estimated weighted regression^{27,43} (LOESS) (**Supplemental Figure 16**). The loess
807 weights are proportional to the product of the inverse variance of each estimated coefficient and
808 the tricube distance function of nearby ages to smooth adjacent ages more closely together
809 proportional to the cube of their distance *d* from the age in question:

810
$$D = \text{abs}(\text{age} - \text{age}_i)$$

811 We consider the neighboring unsmoothed coefficients as those within an adjusted window
812 length, and if the age in question is within 5 years of the minimum or maximum age, we extend
813 the adjusted window by 5 years.

814
$$\text{neighbors} = \text{which}(D_i \leq \text{adjusted}_{\text{window}_{\text{width}}})$$

815
$$\text{weights}_{\text{tricube}} = 1 - \left(\frac{D_{\text{neighbors}}}{\text{window}_{\text{width}}} \right)^3$$

816
$$\text{weights} < - \text{weight}_{\text{tricube}} * \frac{1}{\sigma^2}.$$

817 We then use weighted least square regression to adjust the coefficient as the weighted sum of
818 neighboring coefficients where the design matrix X is the 'N' neighbor' by degree +1 matrix X
819 and y is the N x 1 vector of unsmoothed coefficients.

820
$$WX = \sqrt{\text{weights}} X$$

821
$$Wy = \sqrt{(\text{weights})} * \text{coefficients}[\text{neighbors}]$$

822
$$\beta < -(WX'WX)^{-1}WX'Wy$$

823
$$\text{smoothed}_{\text{coefficients}_i} = \sum \beta + \beta * \text{Age}_{\{i\}} \dots \beta * \text{Age}_i^d$$

824

825 A vignette showing this process on a sample calculation is shown here

826 <https://surbut.github.io/MSGene/vignette.html>. Furthermore, flexible window choices and
827 polynomial degrees can be found here: <https://surbut.shinyapps.io/testapp/>. All analyses were
828 performed with R (version 4.3.1) and our software is written as R code with implementation and
829 vignettes at <https://github.com/surbut/MSGene>.

830

831 **Standard Error of Projection**

832 We bootstrap our training data 50 times and extract the corresponding means and
833 standard errors of each projection across bootstrapping iterations. We compute the remaining
834 lifetime risk setting the maximum age considered as 80 according to the density of observations
835 in our training data, and impute a relative risk (RR) of CAD from statins of 0.20^{30,44,45}; notably,
836 the RR may be larger for some groups, such as those with elevated CAD PRS^{32,46}, and for
837 longer periods of time and thus this reflects a conservative estimate⁴⁷. We apply this benefit only
838 to individuals *not* previously on statins.

839 For the RMSE, we report the standard error of the mean across strata. For proportions,
840 we report the standard error of the sample proportion as $\sqrt{(\hat{p}q/n)}$ where \hat{p} represents the
841 sample proportion.

842 ***Precision and Discrimination analysis***

843 For each age, we compare the average predicted score by genomic (<20%, 20–80%,
844 and >80%) and sex strata, and report the root mean squared error (RMSE) as the difference in
845 the average empirical and predicted cumulative incidence rate for each PRS and sex group as
846 detailed in the Supplementary Methods.

847

$$848 \quad RMSE = \text{sqr}t(\sqrt{\text{Empirical Incidence} - \text{mean}(\text{Predicted Rate}_{PRS \times sex})}.$$

849

850 For the area under the receiver operator curve (AUC-ROC) and precision-recall analysis, we
851 compute the area under each curve using each score as a predictor of cumulative case or
852 control status computed using values for individuals at each year plotted.

853

854 ***States and competing risk***

855 The unique nature of our multistate model features eight mutually exclusive states and restricts
856 one-year transitions as follows (**Fig. 1**), with death as the final absorbing state from which one
857 cannot exit. At any age across the life course, cumulative one-step transitions can be assessed
858 (**Fig. 1**). Possible transitions are as follows:

- 859 1. Health to a single risk factor (Htn, Hld, Dm), CAD or death;
- 860 2. Single risk factor to corresponding double risk factor, CAD or death;
- 861 3. Double risk factor to triple risk factor, CAD or death;
- 862 4. Triple risk factor to CAD or death;
- 863 5. CAD to death.

864

865 ***Predictions with age as the time scale***

866 Our model inferences are made per-year using the individuals who are in a particular risk state
867 at a given age (**Fig. 2, Supp. Fig. 12**). Predictions can, therefore, be made over a requested
868 time interval using the product of age-specific risks for which coefficients were estimated from
869 individuals who were in the at-risk subset during a given period. While enrollment in the UK

870 Biobank required that an individual be alive at age 40 to enroll for genotyping, it did not require
871 that the individual be risk factor-free, and therefore we use this information to assign individuals
872 into risk categories for inference from age 40 onward. We exclude individuals with CAD at
873 baseline from our predictions.

874 *Comparison to 10-year PCE and 30-year Framingham CAD risks*

875 For comparison of time-dependent 10-year risk, we use the 2018 PCE with baseline covariates
876 (total cholesterol, HDL-cholesterol and systolic blood pressure, current smoking) obtained from
877 UKB enrollment data and update each prediction²⁶ with time-varying age, diabetes, and
878 medication use according to available records. This technique was used in the Framingham 30-
879 year risk development to validate new longer window estimates in which age was iteratively
880 updated with all other risk factors at their baseline values.²⁶

881 For comparison of calibration to 30-year risk, we used the 2009 complete (non-BMI)
882 Framingham 30-year equation (FRS30) and update each prediction²⁶ with time-varying age,
883 diabetes, and anti-hypertensive use according to available records, consistent with detailed
884 formulae within the FRS30. Given the differing populations, we recalibrated⁴⁸ according to the
885 mean levels of each covariate and baseline hazard in the UKB sample (FRS30RC). For fair
886 comparison, we report our results against FRS30RC given its improved calibration in our cohort
887 (**Supp. Fig. 17**). Precision and discrimination analysis described as follows. We compute and
888 display the predicted 30-year risk for individuals from ages 40–70 according to this model.

889 *Time-dependent model assessment*

890 We first use the age and state-specific predicted risk scores for each individual - which
891 arise from our MSGene system of smoothed logistic regressions - as covariates in a time-
892 dependent Cox model, in which an individual is featured in non-overlapping intervals with their
893 respective score and event status. In the evaluation stage, we conservatively left censor
894 individuals until enrollment. We also calculate the minimum age at which an individual would
895 exceed pre-specified risk thresholds for MSGene, PCE, and FRS30. We divide every
896 individual's observed trajectory into non-overlapping intervals, indicating when one or all
897 thresholds are achieved and when an event occurs. For example, if an individual is observed
898 from ages 40-70 and exceeds one risk score at age 45 and the other at age 52 and has an
899 event at age 68, his period of study will be divided into 4 intervals: the period from age 40 to 44
900 in which he exceeds the threshold with neither score, the period from 45-51 in which he
901 exceeds the threshold only with score 1, the period from 52 to 67 in which exceeds with both

902 scores, and the period from 68 to 80 in which he has had an event and exceeded in both score.
903 We left censor in this analysis at age of enrollment. We fit independent time-dependent Cox
904 models⁴⁹ to this expanded data set, and again conservatively left censor until enrollment. For
905 both analyses, we report the concordance index (Harrell's-C) with confidence intervals derived
906 from bootstrapping iterations.⁵⁰

907 *Internal and external model assessment*

908 We internally assess the calibration (RMSE) (**Supp. Table 1**) of models using a finite number of
909 covariates for eight state-specific transitions built on a training set and independently assess on
910 our testing set. External validation was performed by comparing the model fits estimated in the
911 UKB with 10-year and lifetime risk estimates from young adults in the Framingham Heart Study
912 Offspring cohort (FOS)⁵¹ (**Supp. Fig. 8**) for whom genetic data are available. This is a
913 community-based Northeastern United States cohort that was recruited in 1971, median age
914 [IQR] 33.0 years [27.0, 41.0] and followed through 2013. Clinical data and incident disease for
915 3836 participants, and genetic data for a subset (2611), were available through the database of
916 Genotypes and Phenotypes (dbGaP; accession phs000007.v33.p14). We compare these with
917 the PCE and FRS30 (original score, calibrated for this population) estimates calculated at Exam
918 1 and compute the RMSE and AUC over the 30-year follow-up period. Informed consent was
919 obtained from all participants, and secondary data analyses of dbGAP based FOS and UKB
920 were approved by the Mass General Brigham Institutional Review Board applications
921 2016P002395 and 2021P002228.

922 **Calculating Net Reclassification**

923
924 For net reclassification indices, at each age of consideration, we defined NRI_{event} as the net
925 proportion of cases correctly reclassified by MSGene Lifetime (MSGene_{LT} >10%) as compared
926 to a ten-year PCE:

927
928
$$NRI_{event} = \frac{MSGene_{LT} > 10\% \cap PCE < 5\% \cap CAD - MSGene_{LT} < 5\% \cap PCE > 5\% \cap CAD}{Develops CAD}$$

929 We defined $NRI_{non-event}$ as the net proportion of controls correctly reclassified by MSGene lifetime
930 risk <10%:

931
$$NRI_{non-event}$$

$$932 \quad \frac{MSGene_{LT} < 10\% \cap PCE > 5\% \cap No CAD - MSGene_{LT} > 10\% \cap PCE < 5\% \cap No CAD}{Does not develop CAD}$$

933 **Marginal Calculation**

934 We also allow, for the absorbing states of CAD and death, the possibility of computing the
935 probability of progressing through any out ('marginal') to CAD. The calculation of progressing to
936 state K from state J through any path over N years is the product of N transition matrices **T** in
937 which the **j,k** element for matrix **T_{ia}** is the probability of progressing from state **j** to **k** at age **a** for
938 individual of covariate profile **i**:

$$940 \quad Marginal\ Interval\ risk = \prod_{A1}^{A2} T_{iajk}$$

941 For every individual, we constrain the row sums to sum to 1 so that the marginal probability
942 across states cannot exceed 1. For absorbing states, the k,k probability is 1. This vignette is
943 available at <https://surbut.github.io/MSGene/usingMarginal.html>.

944

945 **Data Availability**

946 All code for running the MSGene model is available at <https://github.com/surbut/MSGene>.
947 Vignettes for running the analyses are available at
948 <https://surbut.github.io/MSGene/vignette.html> and
949 <https://surbut.github.io/MSGene/usingMarginal.html>. Shiny app for calculating interval risk is
950 available at <https://surbut.shinyapps.io/risk/>. UK Biobank data is available upon application
951 through the UKB Showcase <https://www.ukbiobank.ac.uk>. Framingham Offspring Data is
952 available through dbGap access by investigator application.
953