# Multi-ancestry polygenic risk scores for venous thromboembolism

Yon Ho Jee[1], Florian Thibord[2,3], Alicia Dominguez[4], Corriene Sept[5], Kristin Boulier[6], Vidhya Venkateswaran[7], Yi Ding[6], Tess Cherlin[8], Shefali Setia Verma[8], Valeria Lo Faro[9,10], Traci M. Bartz[11], Anne Boland[12,13], Jennifer A. Brody[14], Jean-Francois Deleuze[12,13,15], Joseph Emmerich[16,17], Marine Germain[18], Andrew D. Johnson[2,3], Charles Kooperberg[19], Pierre-Emmanuel Morange[20], Nathan Pankratz[21], Bruce M. Psaty[22-24], Alexander P. Reiner[23,19], David M. Smadja[25,26], Colleen M. Sitlani[22], Pierre Suchon[20], Weihong Tang[27], David-Alexandre Trégouët[18], Sebastian Zöllner[4], Bogdan Pasaniuc[7], Scott M. Damrauer[28-30], Serena Sanna[31,32], Harold Snieder[9], Lifelines Cohort Study, Christopher Kabrhel[33], Nicholas L. Smith[24,34,35], Peter Kraft[36], INVENT Consortium.


[1]Department of Epidemiology, Harvard T.H. Chan School of Public Health, MA, USA.
[2]Population Sciences Branch, Division of Intramural Research, National Heart, Lung and Blood Institute, MD, USA.
[3]The Framingham Heart Study, 73 Mt. Wayte Ave, Suite #2, Framingham, MA, 01702 USA
[4]Department of Biostatistics, University of Michigan, Ann Arbor, MI, USA.
[5]Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA
[6]Bioinformatics Interdepartmental Program, University of California Los Angeles, Los Angeles, CA, USA.
[7]Department of Oral Biology, University of California Los Angeles School of Dentistry, Los Angeles, CA, USA.
[8]Department of Pathology and Laboratory Medicine, University of Pennsylvania Perelman School of Medicine, Philadelphia, PA, USA.
[9]Department of Epidemiology, University of Groningen, University Medical Center Groningen, Groningen, The Netherlands.
[10]Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala University, Uppsala, Sweden
[11]Cardiovascular Health Research Unit, Departments of Biostatistics and Medicine, University of Washington, 4333 Brooklyn Ave, Seattle, WA 98195
[12]Université Paris-Saclay, CEA, Centre National de Recherche en Génomique Humaine, Evry, France.
[13]Laboratory of Excellence in Medical Genomics, GENMED, Evry, France.
[14]Cardiovascular Health Research Unit, Department of Medicine, University of Washington, 4333 Brooklyn Ave, Seattle, WA 98195
[15]Centre d'Etude du Polymorphisme Humain, Fondation Jean Dausset, Paris, France.
[16]Department of Vascular Medicine, Paris Saint-Joseph Hospital Group, University of Paris, Paris, France.
[17]UMR1153, INSERM CRESS, Paris, France.
[18]University of Bordeaux, INSERM, Bordeaux Population Health Research Center, UMR 1219, Bordeaux, France.
[19]Division of Public Health Sciences, Fred Hutchinbson Cancer Center, Seattle WA 98109
[20]Aix-Marseille University, INSERM, INRAE, Centre de Recherche en CardioVasculaire et Nutrition, Laboratory of Haematology, CRB Assistance Publique – Hôpitaux de Marseille, HemoVasc, Marseille, France
[21]Department of Laboratory Medicine and Pathology, University of Minnesota, Minneapolis, Minnesota, 55455, USA
[22]Cardiovascular Health Research Unit, Department of Medicine, University of Washington, 4333 Brooklyn Ave, Seattle, WA 98195

46   [23]Department of Epidemiology, University of Washington, 4333 Brooklyn Ave, Seattle, WA 98195

47   [24]Department of Health Systems and Population Health, University of Washington, 4333 Brooklyn Ave,
48   Seattle, WA 98195

49   [25]Innovative Therapies in Hemostasis, Université de Paris, INSERM, F-75006 Paris, France.

50   [26]Hematology Department and Biosurgical Research Lab (Carpentier Foundation), Assistance Publique
51   Hôpitaux de Paris, Centre-Université de Paris (APHP-CUP), F-75015 Paris, France

52   [27]Division of Epidemiology and Community Health, School of Public Health, University of Minnesota,
53   Minneapolis, Minnesota, 55454, USA

54   [28]Department of Genetics, University of Pennsylvania Perelman School of Medicine, Philadelphia, PA, USA.

55   [29]Department of Surgery, Department of Genetics, and Cardiovascular Institute, Perelman School of
56   Medicine, University of Pennsylvania, Philadelphia PA

57   [30]Department of Surgery, Corporal Michael Crescenz VA Medical Center, Philadelphia PA

58   [31]University of Groningen, UMCG, Department of Genetics, Groningen, the Netherlands

59   [32]Institute for Genetics and Biomedical Research, National Research Council, Monserrato, Italy

60   [33]Center for Vascular Emergencies, Department of Emergency Medicine, Massachusetts General Hospital,
61   Harvard Medical School, Boston, MA, USA.

62   [34]Kaiser Permanente Washington Health Research Institute, Kaiser Permanente Washington, Seattle WA
63   98101, USA

64   [35]Seattle Epidemiologic Research and Information Center, Department of Veterans Affairs Office of
65   Research and Development, Seattle WA 98108, USA

66   [36]Transdivisional Research Program, Division of Cancer Epidemiology and Genetics, National Cancer
67   Institute, National Institutes of Health, MD, USA.

68

69   **Contact**: Peter Kraft, Transdivisional Research Program, Division of Cancer Epidemiology and
70   Genetics, National Cancer Institute, National Institutes of Health, MD, USA.
71   (phillip.kraft@nih.gov)

## Abstract

Venous thromboembolism (VTE) is a significant contributor to morbidity and mortality, with large disparities in incidence rates between Black and White Americans. Polygenic risk scores (PRSs) limited to variants discovered in genome-wide association studies in European-ancestry samples can identify European-ancestry individuals at high risk of VTE. However, there is limited evidence on whether high-dimensional PRS constructed using more sophisticated methods and more diverse training data can enhance the predictive ability and their utility across diverse populations. We developed PRSs for VTE using summary statistics from the International Network against Venous Thrombosis (INVENT) consortium GWAS meta-analyses of European- (71,771 cases and 1,059,740 controls) and African-ancestry samples (7,482 cases and 129,975 controls). We used LDpred2 and PRSCSx to construct ancestry-specific and multi-ancestry PRSs and evaluated their performance in an independent European- (6,261 cases and 88,238 controls) and African-ancestry sample (1,385 cases and 12,569 controls). Multi-ancestry PRSs with weights tuned in European- and African-ancestry samples, respectively, outperformed ancestry-specific PRSs in European- (PRSCSX$_{EUR}$: AUC=0.61 (0.60, 0.61), PRSCSX_combined$_{EUR}$: AUC=0.61 (0.60, 0.62)) and African-ancestry test samples (PRSCSX$_{AFR}$: AUC=0.58 (0.57, 0.6), PRSCSX_combined $_{AFR}$: AUC=0.59 (0.57, 0.60)). The highest fifth percentile of the best-performing PRS was associated with 1.9-fold and 1.68-fold increased risk for VTE among European- and African-ancestry subjects, respectively, relative to those in the middle stratum. These findings suggest that the multi-ancestry PRS may be used to identify individuals at highest risk for VTE and provide guidance for the most effective treatment strategy across diverse populations.

## Introduction

94

95 Venous thromboembolism (VTE) is among the top five most common vascular diseases in most

96 countries (1). The estimated lifetime risk of VTE is 8% among US adults (2). Approximately 20%

97 of individuals die within 1 year of a VTE diagnosis often from the provoking conditions, and

98 complications are common among survivors (3). Thus, the development of tools that stratify

99 people according to their risk of developing VTE is helpful, which could inform risk-stratified

100 prevention strategies that contribute to reducing the burden of VTE.

101 Polygenic risk scores (PRS) are useful tools for approximating the cumulative genetic

102 susceptibility to complex traits and diseases. PRSs based on the independent genome-wide

103 significant variants discovered in genome-wide association studies (GWAS) European-ancestry

104 samples (4–9) have been demonstrated to identify individuals at high risk of VTE (10,11).

105 However, there is limited evidence on whether high-dimensional PRS that are not restricted to

106 genome-wide significant variants can enhance the predictive ability.

107 In the USA, the incidence of VTE is approximately 65% higher in those who identify as Black

108 Americans than White Americans (12,13). Polygenic risk prediction models for VTE could be

109 particularly important among Black Americans, as a clinical tool to reduce this disparity in VTE

110 risk. (This does not preclude research into structural inequities and social determinants of

111 health, which might inform policy interventions to reduce disparities between Black and White

112 Americans.) However, previously developed VTE PRS have been optimized for European-

113 ancestry populations, and their utility in other populations is unknown. In particular, we are

114 unaware of any efforts to develop VTE PRS specifically for Black Americans.

115    We developed ancestry-specific and multi-ancestry PRSs for VTE leveraging large GWAS meta-

116    analyses in European-and African-ancestry samples. We validated these PRSs by estimating

117    relative VTE risks across PRS quintiles in five independent U.S.-based studies. We focus on PRS

118    including common variants (minor allele frequencies above 1%) due to difficulties measuring or

119    imputing low frequency or rare variants from GWAS data or imprecision of estimating rare

120    variant associations. Thus our PRSs complement known low frequency variants (such as rs6205

121    in *F5*) or known clinical and behavioral risk factors. Here we concentrate on developing PRSs

122    that perform well in diverse populations. Future work will be needed to (a) develop and

123    evaluate models that combine these PRSs with low-frequency and rare variants and other risk

124    factors and (b) assess the clinical utility of VTE risk models for targeted prevention, screening,

125    or treatment (14,15).

126

## Results

## Study sample

129    The overall study design is illustrated in **Figure 1**. Our PRS development consisted of two steps:

130    training ancestry-specific PRS and tuning multi-ancestry PRS. We trained ancestry-specific PRSs

131    using European- and African ancestry GWAS summary statistics from the INVENT consortium

132    and two Bayesian methods (LDPRED2(14) and PRSCSx(15)). We then tuned the constructed

133    multi-ancestry PRSs by regressing VTE case-control status on a linear combination of the two

134    ancestry-specific PRSs in two separate tuning samples: one European-ancestry tuning sample

135    (1,329 cases and 1,324 controls) and one African-ancestry tuning sample (238 cases and 3,589

136     controls). The testing data set comprised 6,781 cases and 103,016 controls of European

137     ancestry and 1,385 cases and 12,569 controls of African ancestry from five independent studies.

138     Table S1 presents a brief summary of participating studies and biobanks, including basic

139     information about each study or biobank (location, institute, cohort size, and sample recruiting

140     approach), participants (ancestry and age), and genotypes (genotyping platforms and

141     imputation reference).

142     **Comparing PRS distributions across populations**

143     Four single-ancestry PRSs and four multi-ancestry PRSs for VTE were constructed using LDpred2

144     and PRSCSx and validated in independent European ancestry and African ancestry  individuals:

145     (i) LDpred2 trained using European-ancestry GWAS summary statistics (LDpred2$_{EUR}$); (ii)

146     LDpred2 trained using African-ancestry summary statistics (LDpred2$_{AFR}$); (iii) PRS-CS trained

147     using European-ancestry summary statistics (PRSCSX $_{EUR}$); (iv) PRSCS trained using African

148     ancestry summary statistics (PRSCSX $_{AFR}$); and (v) LDpred2$_{EUR}$ + LDpred2$_{AFR}$ with weights tuned in

149     an independent European-ancestry tuning sample; (vi) LDpred2$_{EUR}$ + LDpred2$_{AFR}$ with weights

150     tuned in and independent African-ancestry tuning sample (LDpred2_combined$_{AFR}$); (vii) PRSCSX

151     $_{EUR}$ + PRSCSX $_{AFR}$ with weights tuned in the European-ancestry tuning sample (PRSCSX_combined

152     $_{EUR}$); (viii) PRSCSX $_{EUR}$ + PRSCSX $_{AFR}$ with weights tuned in the African-ancestry tuning sample

153     (PRSCSX_combined$_{AFR}$). All PRSs had higher means in cases than controls in the test data sets

154     (**Table 1**). Among the European-ancestry VTE cases, the mean PRS was higher for the PRS tuned

155     in European-ancestry samples than for the PRS tuned in African-ancestry samples. The

156     difference was higher for the ancestry-specific PRS (LDpred2$_{EUR}$: 0.39 vs LDpred2$_{AFR}$: 0.07,

157     PRSCSX$_{EUR}$: 0.42 vs PRSCSX$_{AFR}$: 0.31) than for the multi-ancestry PRS (LDpred2_combined $_{EUR}$:

158  0.39 vs Dpred2_combined $_{AFR}$: 0.38, PRSCSX_combined $_{EUR}$: 0.44 vs PRSCSX_combined $_{AFR}$: 0.41).

159  Similarly, among the African-ancestry VTE cases, the mean PRS was higher for the African-

160  ancestry-tuned PRS than for the European-ancestry-tuned PRS, with larger difference for the

161  population-specific PRS (LDpred2$_{EUR}$: 0.18 vs Dpred2$_{AFR}$: 0.19, PRSCSX$_{EUR}$: 0.22 vs PRSCSX$_{AFR}$:

162  0.28) than the multi-ancestry PRS (LDpred2_combined $_{EUR}$: 0.19 vs Dpred2_combined $_{AFR}$: 0.23,

163  PRSCSX_combined $_{EUR}$: 0.26 vs PRSCSX_combined $_{AFR}$: 0.30).

164  **Evaluation of PRS and VTE risk across populations**

165  **Table 2** shows the estimated OR per SD increase of PRS and AUC for VTE in the test set

166  individuals of European- and African ancestry. For the ancestry-specific PRS, LDpred2$_{EUR}$ and

167  LDpred2$_{AFR}$ were constructed using 604,741 SNPs and 1,184,805 SNPs, respectively, and

168  PRSCSX$_{EUR}$ and PRSCSX$_{AFR}$ were constructed using 591,788 SNPs and 586,660 SNPs, respectively.

169  Multi-ancestry PRS were developed as a linear combination of the ancestry-specific PRS,

170  resulting in 1,212,566 SNPs for LDpred2 and 598,977 SNPs for PRSCSX. The multi-ancestry PRSs

171  outperformed ancestry-specific PRSs in both European- and African-Ancestry test samples and

172  across training methods (LDpred2, PRSCSx) (**Figure 2, S.Figure 1**). In the European-ancestry test

173  set, multi-ancestry PRS in which the weights were tuned in European ancestry samples

174  performed the best (PRSCSX_combined $_{EUR}$: AUC=0.61 (0.6, 0.62), OR=1.48 (1.45, 1.52),

175  LDpred2_combined $_{EUR}$: AUC=0.60 (0.59, 0.61), OR=1.42 (1.39, 1.46)). Similarly, in the African-

176  ancestry test set, a multi-ancestry PRS in which the weights were tuned in African-Ancestry

177  samples performed the best (PRSCSX_combined $_{AFR}$: AUC=0.59 (0.57, 0.60), OR=1.38 (1.30, 1.45);

178  LDpred2_combined $_{AFR}$: AUC=0.57 (0.55, 0.58), OR=1.26 (1.20, 1.33)).

179     The association between the PRSs and VTE risk by PRS percentile are shown in **Figure 3**. The

180     association between the highest fifth percentile of $PRSCSX_{EUR}$ (RR=1.89) and $LDpred2_{EUR}$

181     (RR=1.79) and VTE risk was greater than that of genome-wide significant PRS (RR=1.78). The

182     highest fifth percentile of the best-performing PRS ($PRSCSX\_combined_{EUR}$) was associated with

183     1.9-fold increased risk for VTE among European ancestry subjects compared to the middle

184     stratum (40–50%). Among the African-ancestry samples, the corresponding risk was about 1.68-

185     fold ($PRSCSX\_combined_{AFR}$), which is smaller than that in European ancestry samples.

186     **Inclusion of known low frequency alleles**

187     When we reconstructed PRS including the five genome-wide significant variants, the new PRS

188     performed worse than our original PRS without the five SNPs in European- ($PRSCSX\_combined$

189     $_{EUR}$: AUC= 0.57 (0.56, 0.59), $LDpred2\_combined_{EUR}$: AUC= 0.52 (0.50, 0.53)) and in African-

190     ancestry test samples ($PRSCSX\_combined_{AFR}$: AUC= 0.59 (0.58, 0.60), $LDpred2\_combined_{AFR}$:

191     AUC= 0.56 (0.55, 0.57)) (**S.Figure 2**). This is likely because the five SNPs are rare in one or both

192     populations (average MAF in European ancestry=0.1, African ancestry=0), and our tuning

193     samples are small, resulting in noisy weights. Future studies with larger and more diverse

194     training samples and further tuning steps are needed to learn better multi-ancestry PRS weights.

195

196     **Discussion**

197     Multi-ancestry PRSs outperformed population specific PRSs in U.S. European- and African-

198     ancestry samples, with a greater improvement in African-ancestry samples. The highest fifth

199     percentile of the best performing multi-ancestry PRS in the European ancestry test samples was

8

200    associated with an approximately 2-fold increased risk for VTE relative to the middle stratum

201    among European-ancestry subjects. The corresponding risk was smaller (1.7-fold) among the

202    African-ancestry subjects, but still non-negligible and higher than any single-ancestry PRS,

203    highlighting that multi-ancestry PRS may be used to identify individuals at highest risk for VTE

204    events. These data may also be useful in guiding primary prevention and treatment strategies

205    across populations, although we stress that demonstrating PRS discrimination is not sufficient

206    to establish clinical utility, which requires consideration of risks and benefits of specific

207    proposed interventions (14,15).

208

209    To our knowledge, this is the first attempt to develop PRS of VTE specific to African-ancestry

210    populations. Clinical evaluation of PRS is needed in African-ancestry populations, where the

211    burden of VTE is growing due to its increase in VTE incidence. Our PRS, developed and validated

212    in African-ancestry samples, could be a step towards risk-based clinical management of VTE

213    among Black Americans, as a complement to primary prevention efforts. Black Americans and

214    other population groups suffer social disadvantage and lifestyle risk factors that could be a

215    strong contributors to the disparities in VTE (16). Encouragingly, healthy lifestyle factors were

216    associated with a lower incidence of VTE among people at high genetic risk for VTE (17). Hence,

217    as with most diseases, primary prevention efforts directed at lifestyle interventions to reduce

218    weight or increase activity would have the great potential to reduce the societal burden of VTE.

219    Further research should determine best approaches to VTE prevention that improve health

220    equity.

221

222    A recent GWAS meta-analysis demonstrated that European-ancestry individuals at or above the

223    top fifth percentile of a PRS comprised of 37 genome-wide significant variants had a 3.2-fold

224    greater risk for VTE (OR: 3.19; 95% CI: 2.89-3.52) relative to half of the population in the middle

225    of the range (8). More recently, a PRS using the 100 lead variants identified in a larger European

226    ancestry meta-analysis showed AUC=0.620 (95% CI, 0.616–0.625) (9). Since these previous PRS

227    include low MAF variants with large effect sizes (e.g., rs6025: transancestry OR=2.39 (8) on *F5*

228    gene), the performance of these previous PRSs and our PRSs is not directly comparable. It is

229    worth noting that our PRS was built using genome-wide common variants and was designed to

230    be transportable between European- and African-ancestry individuals, which can be useful for

231    settings with diverse genetic background. The PRSs presented here complement the low-

232    frequency, large-effect variants and clinical and behavioral risk factors; future work should

233    develop and evaluate comprehensive risk models combining multi-ancestry PRS, low-frequency

234    variants and other risk factors.

235    The major strength of the study is that it is the first attempt to develop and validate multi-

236    ancestry PRS for VTE, providing potential utility of PRS in VTE prevention among African-

237    ancestry populations, where the VTE burden is high. In addition, we validated the PRS in the

238    five independent biobanks from GBMI using harmonized analysis framework (e.g. phenotype

239    definitions, ancestry assignments, and PRS construction).

240    There are several limitations in our study. First, we have focused on common SNPs, specifically

241    HapMap3 SNPs for VTE PRS construction. As a result, information from rarer variants missing in

242    the LD reference panel may not be captured in other non-European ancestries. Second, the

243    lower predictive ability of VTE PRS in African-ancestry samples can be explained by smaller

244    sample size of African-ancestry VTE meta-analysis GWAS, which is 10 times smaller than

245    European GWAS. Third, there remains a multitude of factors that may contribute to cross-

246    biobank heterogeneity including phenotype precision, cohort-level disease prevalence, and

247    environmental factors. We have provided analysis results by cohort (**Supplementary Figure 1**).

248    **Conclusions**

249    We found that multi-ancestry PRS for VTE outperformed population-specific PRS, especially in

250    African ancestry populations with relatively small GWAS sample sizes. These findings suggest

251    that the multi-ancestry PRS may be used to identify individuals at highest risk for VTE event and

252    provide guidance for the most effective treatment strategy across populations.

253

254

255    **Materials and Methods**

256    **Study populations**

257    We trained the PRS using summary statistics from the International Network against Venous

258    Thrombosis (INVENT) consortium cross-ancestry GWAS meta-analyses of European- (71,771

259    VTE cases and 1,059,740 controls) and African-ancestry samples (7,482 VTE cases and 129,975

260    controls) (9). The meta-analysis is based on prospective cohorts and case-control data from 30

261    studies.

262    Tuning (1,329 cases and 1,324 controls of European-ancestry and 238 cases and 3,589 controls

263    of African-ancestry) and validation data (6,781 cases and 103,016 controls of European ancestry

264    and 1,385 cases and 12,569 controls of African ancestry) came from Nurses' Health Study [NHS]

265    and Health Professional Follow-up Study [HPFS] and 4 Global Biobank Meta-analysis Initiative

266    (GBMI) biobanks (Michigan Genomics Initiative [MGI], UCLA Precision Health Biobank [UCLA],

267    Penn Medicine Biobank [PMBB], and Lifelines) with representation across African and

268    European-ancestry populations included (**Figure 1**). These tuning and validation data were not

269    included in the GWAS used in the training step. The definitions of African- and European-

270    ancestry populations in each study are provided in the **Supplementary Materials**; these

271    definitions typically involve both self-reported race and ethnicity and genetic similarity to a set

272    of (study-specific) labeled reference samples.

273    **Supplementary Table 1** summarizes the study design, genotyping arrays, and the sample size in

274    each study. All studies were approved by the relevant institutional ethics committees and

275    review boards, and all participants provided written informed consent.

276

277    **Statistical methods**

278    **PRS training and tuning**

279    *PRS training and tuning using LDpred2.* We ran LDpred2-auto(14) to construct PRS on HapMap3

280    variants using the INVENT GWAS meta-analysis summary statistics corresponding to each

281    population. We constructed linkage disequilibrium (LD) reference panels for the development

282    of the European-ancestry PRS (LDpred2$_{EUR}$) and African-Ancestry PRS (LDpred2$_{AFR}$) using the

283    EUR and AFR supersamples from the 1000 Genomes Project (Phase 3), respectively.(18) These

284    population-specific PRSs were then linearly combined to construct multi-ancestry PRS

12

285   (LDpred2$_{EUR}$ + LDpred2$_{AFR}$) in which the relative contribution of each PRS was estimated by

286   logistic regression in the tuning dataset of European-ancestry samples (LDpred2_combined$_{EUR}$)

287   and African-ancestry samples (LDpred2_combined$_{AFR}$). Analyses were run using R; code is

288   available at https://github.com/yonhojee/VTE_PRS.

289   *PRS training and tuning using PRSCSx.* We separately applied PRSCSx(15) to the summary

290   statistics from the European- and African-ancestry INVENT VTE GWAS, using the EUR and AFR

291   LD reference panels from the 1000 Genomes Project (Phase 3). The global shrinkage parameter

292   was learnt from the data using a fully Bayesian approach. Ancestry-specific PRSs generated

293   using European- (PRSCSx$_{EUR}$) and African-specific posterior weights (hereafter denoted as

294   PRSCSx$_{AFR}$) were linearly combined to construct multi-ancestry PRS (PRSCSx$_{EUR}$ + PRSCSx$_{AFR}$). The

295   regression coefficients for the linear combination were obtained by fitting a logistic regression

296   model in the tuning data set of European ancestry samples (PRSCSx_combined$_{EUR}$) and African

297   American samples (PRSCSx_combined$_{AFR}$). Analyses were run using Python; code is available at

298   https://github.com/yonhojee/VTE_PRS.

299

300   **PRS validation**

301   In each test dataset, population-specific PRSs were calculated as $PRS_{EUR_i} = \sum \beta_k x_{ik}$ and

302   $PRS_{AFR_i} = \sum \alpha_k x_{ik}$, where $x_{ik}$ is the dosage of risk allele (0-2) at genetic variant $k$ for subject $i$,

303   and $\beta_k$ and $\alpha_k$ are the corresponding weight in European and African PRS, respectively. The

304   estimates of $\beta_k$ and $\alpha_k$ were trained using summary statistics from the INVENT consortium and

305   LDpred2 and PRSCSx as described above.

13

306     We calculated the multi-ancestry PRSs as the linear combination of European- and African-

307     ancestry specific PRS:

$$PRS\_combined_{EUR_i} = \gamma_{AFR\_EUR}PRS_{AFR_i} + \gamma_{EUR\_EUR}PRS_{EUR_i}$$

$$PRS\_combined_{AFR_i} = \delta_{AFR\_AFR}PRS_{AFR_i} + \delta_{EUR\_AFR}PRS_{EUR_i}$$

308     where $PRS_{AFR}$ and $PRS_{EUR}$ are the PRSs trained in single-ancestry GWAS and the $\gamma$ and $\delta$ are

309     "meta-weights" tuned in European- and African-ancestry samples, respectively. SNPs with

310     imputation $R^2 > 0.9$ in training dataset were retained for subsequent analyses. The lists of SNPs

311     and the weights for the PRS computation are available at

312     https://github.com/yonhojee/VTE_PRS.

313     PRSs were standardized within each validation sample to have unit SD in the control subjects.

314     Logistic regression, adjusting for ten principal components and sex, was used to estimate odds

315     ratios (ORs) for association between the standardized PRSs and VTE risk in each testing set. The

316     discrimination of PRS was assessed using area under the receiver operating curve (AUC). The OR

317     per SD and AUC were obtained individually for each study and combined separately for

318     European- and African-ancestry samples using a fixed-effect meta-analysis.

319     All statistical analyses were conducted using R v.4.3.0. Logistic regression and AUC were done

320     using *glm()* and *roc()* in R.

321

322     **The distribution of relative risk of VTE by PRS across populations.**

323    We simulated 100,000 individuals with PRS distribution of N(0,1) multiplied by log OR per SD

324    estimates for each PRS. The simulated PRS was then exponentiated to estimate relative risk

325    estimates and split into the percentile categories: [0–1%] (1-5%], (5-10%], (10–20%], (20–30%],

326    (30–40%], (40–50%] (reference group), (50–60%], (60–70%], (70–80%], (80–90%], (90–95%],

327    (95–99%] and (99–100%].

328

329    **Sensitivity analysis of including known low frequency alleles**

330    Out of the 37 genome-wide significant variants, our current PRSs do not include five variants

331    (rs6025, rs145470028, rs1799963, rs6048, and rs143478537), which would have been filtered

332    out of our analyses for one reason or another (e.g., on the X chromosome, low minor allele

333    frequency [MAF]). These variants are important to be considered in VTE PRS given their large

334    effect sizes (e.g., rs6025: transancestry OR=2.39(8) on *F5* gene). As a sensitivity analysis, we

335    constructed new PRSs, which additionally include these previously reported variants that are i)

336    not included in our PRS due to the low frequency and ii) not in LD with the variants already

337    included in our PRS. The final PRSs were obtained by the linear combination of the original PRS

338    (constructed using common variants only) and the additional SNPs where the coefficients for

339    the original PRS and the additional SNPs were tuned in the independent ancestry-specific

340    samples (See more details in the **Supplementary Materials**).

341

## Conflict of Interest Statement

393     B.M.P. serves on the Steering Committee of the Yale Open Data Access Project funded by
394     Johnson & JOhnson.

395     S.M.D. receives research support from RenalytixAI and Novo Nordisk, outside the scope of the
396     current research. SMD is named as a co-inventor on a Government-owned US Patent
397     application related to the use of genetic risk prediction for venous thromboembolic disease
398     filed by the US Department of Veterans Affairs in accordance with Federal regulatory
399     requirements.

400

401

# References

1. Wendelboe, A.M. and Raskob, G.E. (2016) Global Burden of Thrombosis: Epidemiologic Aspects. *Circ Res*, **118**, 1340–1347.
2. Bell, E.J., Lutsey, P.L., Basu, S., Cushman, M., Heckbert, S.R., Lloyd-Jones, D.M. and Folsom, A.R. (2016) Lifetime Risk of Venous Thromboembolism in Two Cohort Studies. *The American Journal of Medicine*, **129**, 339.e19-339.e26.
3. Søgaard, K.K., Schmidt, M., Pedersen, L., Horváth-Puhó, E. and Sørensen, H.T. (2014) 30-year mortality after venous thromboembolism: a population-based cohort study. *Circulation*, **130**, 829–836.
4. Heit, J.A., Armasu, S.M., Asmann, Y.W., Cunningham, J.M., Matsumoto, M.E., Petterson, T.M. and De Andrade, M. (2012) A genome-wide association study of venous thromboembolism identifies risk variants in chromosomes 1q24.2 and 9q. *J Thromb Haemost*, **10**, 1521–1531.
5. Tang, W., Teichert, M., Chasman, D.I., Heit, J.A., Morange, P.-E., Li, G., Pankratz, N., Leebeek, F.W., Paré, G., de Andrade, M., *et al.* (2013) A Genome-Wide Association Study for Venous Thromboembolism: The Extended Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) Consortium. *Genetic Epidemiology*, **37**, 512–521.
6. Germain, M., Chasman, D.I., de Haan, H., Tang, W., Lindström, S., Weng, L.-C., de Andrade, M., de Visser, M.C.H., Wiggins, K.L., Suchon, P., *et al.* (2015) Meta-analysis of 65,734 individuals identifies TSPAN15 and SLC44A2 as two susceptibility loci for venous thromboembolism. *Am J Hum Genet*, **96**, 532–542.
7. Klarin, D., Busenkell, E., Judy, R., Lynch, J., Levin, M., Haessler, J., Aragam, K., Chaffin, M., Haas, M., Lindström, S., *et al.* (2019) Genome-wide association analysis of venous thromboembolism identifies new risk loci and genetic overlap with arterial vascular disease. *Nat Genet*, **51**, 1574–1579.
8. Lindström, S., Wang, L., Smith, E.N., Gordon, W., van Hylckama Vlieg, A., de Andrade, M., Brody, J.A., Pattee, J.W., Haessler, J., Brumpton, B.M., *et al.* (2019) Genomic and transcriptomic association studies identify 16 novel susceptibility loci for venous thromboembolism. *Blood*, **134**, 1645–1657.
9. Thibord, F., Klarin, D., Brody, J.A., Chen, M.-H., Levin, M.G., Chasman, D.I., Goode, E.L., Hveem, K., Teder-Laving, M., Martinez-Perez, A., *et al.* Cross-Ancestry Investigation of Venous Thromboembolism Genomic Predictors. *Circulation*, **0**, 10.1161/CIRCULATIONAHA.122.059675.
10. D, K., E, B., R, J., J, L., M, L., J, H., K, A., M, C., M, H., S, L., *et al.* (2019) Genome-wide association analysis of venous thromboembolism identifies new risk loci and genetic overlap with arterial vascular disease. *Nat Genet*, **51**, 1574–1579.
11. Kolin, D.A., Kulm, S. and Elemento, O. (2021) Prediction of primary venous thromboembolism based on clinical and genetic factors within the U.K. Biobank. *Sci Rep*, **11**, 21340.
12. Silverstein, M.D., Heit, J.A., Mohr, D.N., Petterson, T.M., O'Fallon, W.M. and Melton, L.J. (1998) Trends in the incidence of deep vein thrombosis and pulmonary embolism: a 25-year population-based study. *Arch Intern Med*, **158**, 585–593.

444    13. Zakai, N.A. and McClure, L.A. (2011) Racial differences in venous thromboembolism. *J*
445        *Thromb Haemost*, **9**, 1877–1882.
446    14. Privé, F., Arbel, J. and Vilhjálmsson, B.J. (2020) LDpred2: better, faster, stronger.
447        *Bioinformatics*, **36**, 5424–5431.
448    15. Ruan, Y., Lin, Y.-F., Feng, Y.-C.A., Chen, C.-Y., Lam, M., Guo, Z., Initiatives, S.G.A., He, L.,
449        Sawa, A., Martin, A.R., *et al.* (2021) Improving Polygenic Prediction in Ancestrally Diverse
450        Populations. Improving Polygenic Prediction in Ancestrally Diverse Populations **(2021)**,
451        2020.12.27.20248738.
452    16. Folsom, A.R., Basu, S., Hong, C.-P., Heckbert, S.R., Lutsey, P.L., Rosamond, W.D. and
453        Cushman, M. (2019) Reasons for Differences in the Incidence of Venous Thromboembolism
454        in Black Versus White Americans. *Am J Med*, **132**, 970–976.
455    17. Evans, C.R., Hong, C.-P., Folsom, A.R., Heckbert, S.R., Smith, N.L., Wiggins, K., Lutsey, P.L.
456        and Cushman, M. (2020) Lifestyle Moderates Genetic Risk of Venous Thromboembolism:
457        The ARIC Study. *Arterioscler Thromb Vasc Biol*, **40**, 2756–2763.
458    18. Auton, A., Abecasis, G.R., Altshuler, D.M., Durbin, R.M., Abecasis, G.R., Bentley, D.R.,
459        Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., *et al.* (2015) A global reference for
460        human genetic variation. *Nature*, **526**, 68–74.
461

## Legends to Figures

**Figure 1.** Overview of development and validation of population-specific and multi-ancestry PRS for venous thromboembolism.

**Figure 2.** AUC and OR for population-specific and multiancestry PRS across populations.

**Figure 3.** Distribution of relative risk of VTE by PRS across populations.

## Legends to Tables

**Table 1.** Mean and standard deviation of standardized polygenic risk scores with VTE risk in the test set individuals of European and African ancestry.

**Table 2.** Association of polygenic risk scores and VTE risk in the test set individuals of European and African ancestry.

473 **Figure 1.** Overview of development and validation of population-specific and multi-ancestry PRS for venous thromboembolism.



474

475 PRS development consisted of two steps: training ancestry-specific PRS and tuning multi-ancestry PRS. We trained ancestry-specific
476 PRSs using European- and African ancestry GWAS summary statistics from the INVENT consortium and two Bayesian methods
477 (LDPRED2 and PRSCSx). We then tuned the constructed multi-ancestry PRSs by regressing VTE case-control status on a linear
478 combination of the two ancestry-specific PRSs in two separate tuning samples: one European-ancestry tuning sample and one
479 African-ancestry tuning sample. NHS, Nurses' Health Study; HPFS, Health Professional Follow-up Study; MGI, Michigan Genomics
480 Initiative; UCLA, UCLA Precision Health Biobank; PMBB, Penn Medicine Biobank.

481    **Figure 2.** AUC and OR for population-specific and multiancestry PRS across populations.



482

483

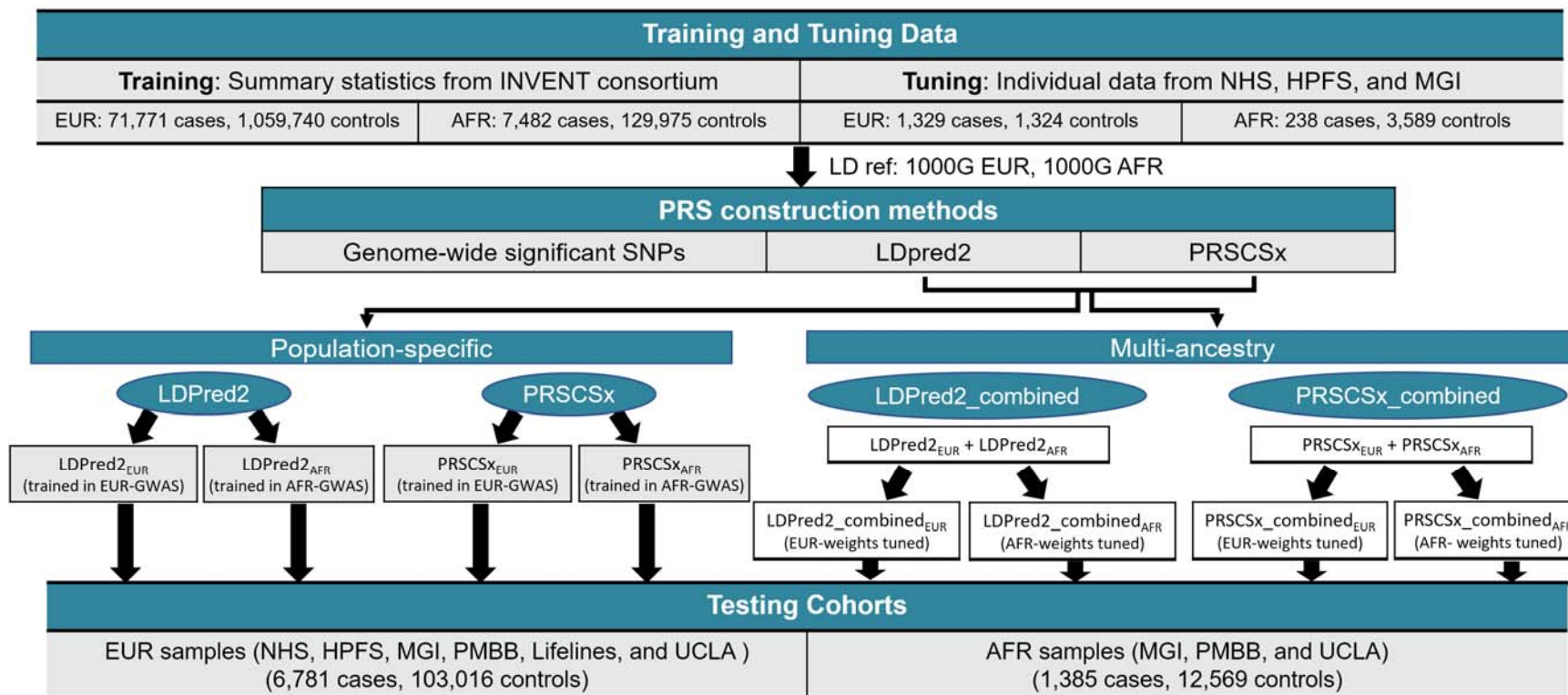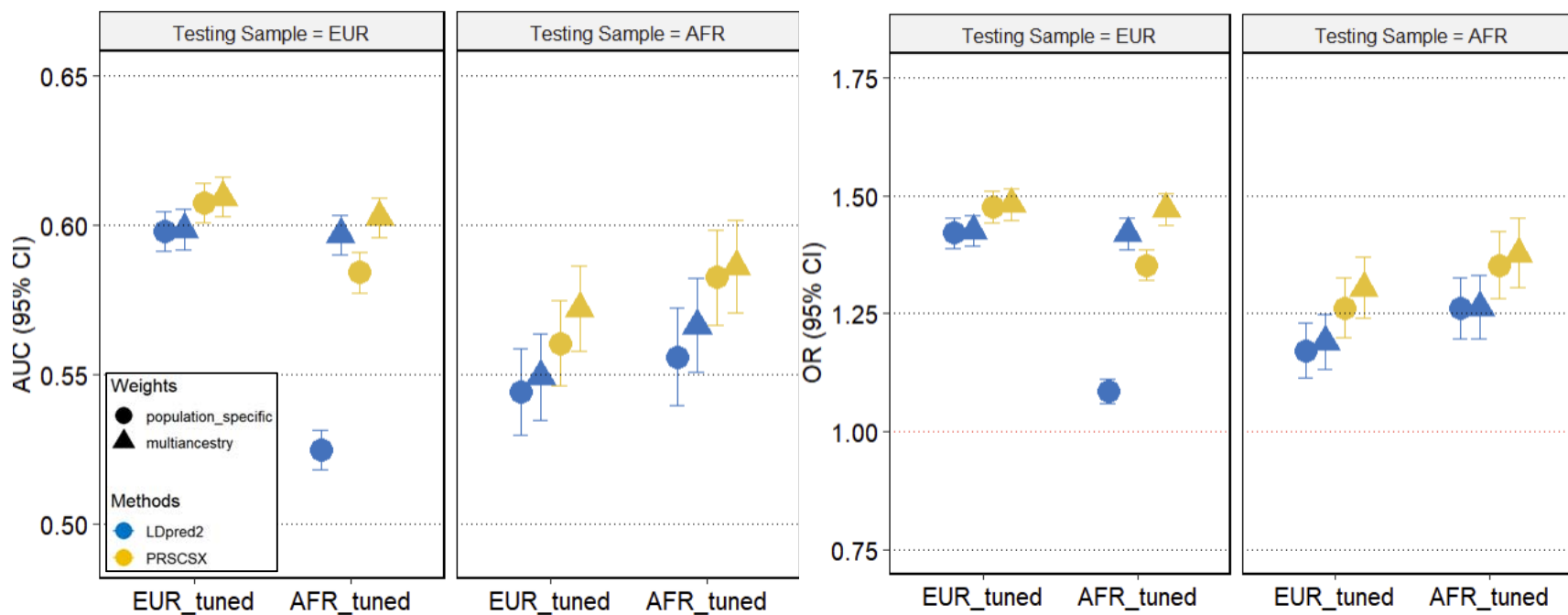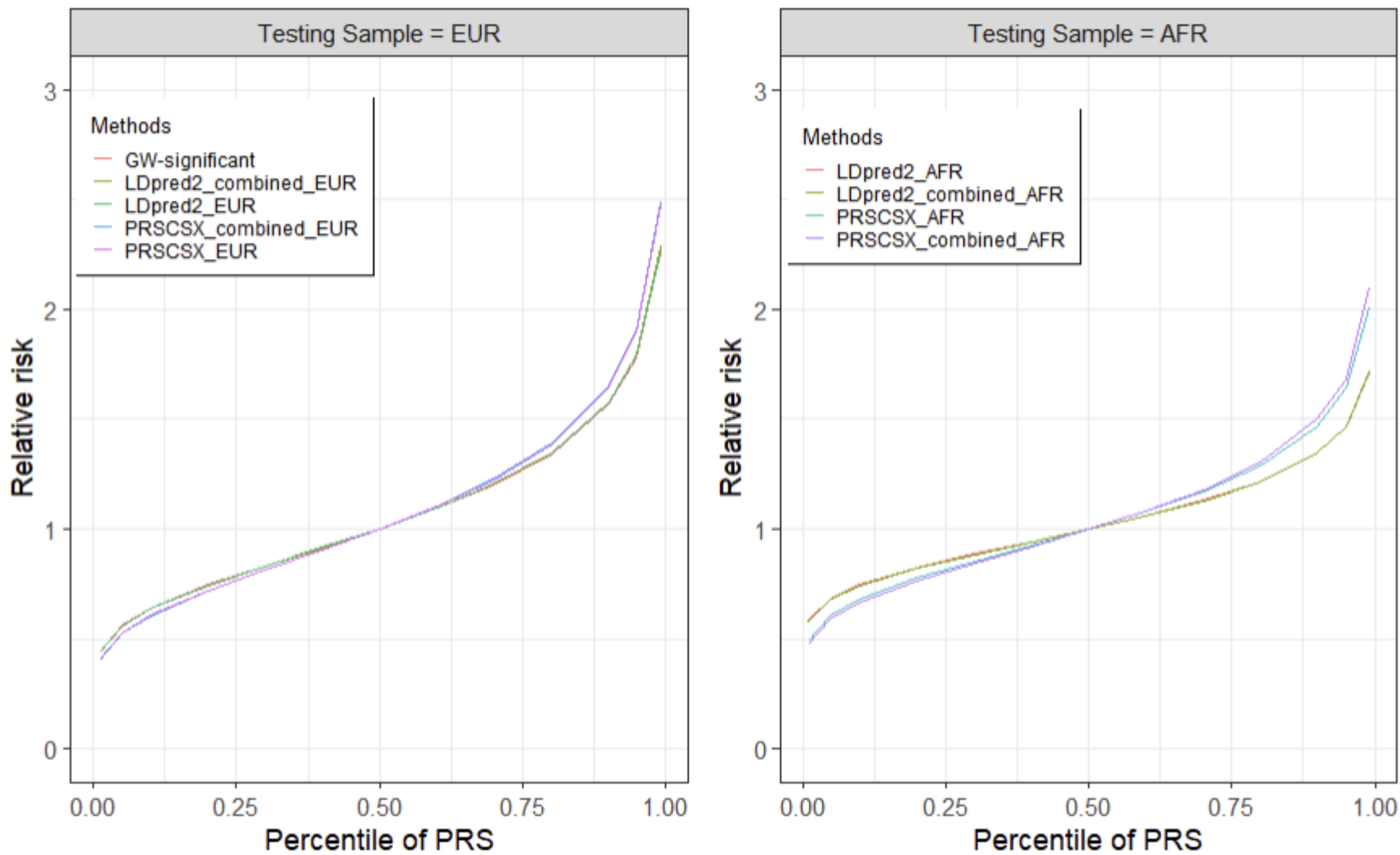484    **Figure 3.** Distribution of relative risk of VTE by PRS across populations.



485

486

487    **Table 1.** Mean and standard deviation of standardized polygenic risk scores with VTE risk in the test set individuals of European and
488    African ancestry.

| | European | | African | |
|---|---|---|---|---|
| | Cases (n=6,781) | Control (n=103,016) | Cases (n=1,385) | Control (n=12,569) |
| Mean (SD) of age at recruitment, in years | 56.9 (13.3) | 52.1 (14.4) | 56.5 (14.6) | 50.8 (16.2) |
| Mean (SD) of LDpred2$_{EUR}$ | 0.39 (1.07) | 0 (1) | 0.18 (1.02) | 0 (1) |
| Mean (SD) of LDpred2$_{AFR}$ | 0.07 (1) | 0 (1) | 0.19 (1.11) | 0 (1) |
| Mean (SD) of PRSCSX$_{EUR}$ | 0.42 (1.07) | 0 (1) | 0.22 (1.03) | 0 (1) |
| Mean (SD) of PRSCSX$_{AFR}$ | 0.31 (1.03) | 0 (1) | 0.28 (1.11) | 0 (1) |
| Mean (SD) of LDpred2_combined$_{EUR}$ | 0.39 (1.07) | -0.02 (1) | 0.19 (1.02) | 0 (1) |
| Mean (SD) of LDpred2_combined$_{AFR}$ | 0.38 (1.06) | 0 (1) | 0.23 (1.04) | 0 (1) |
| Mean (SD) of PRSCSX_combined$_{EUR}$ | 0.44 (1.19) | 0 (1) | 0.26 (1.06) | 0 (1) |
| Mean (SD) of PRSCSX_combined$_{AFR}$ | 0.41 (1.07) | 0 (1) | 0.3 (1.09) | 0 (1) |

489    SD, standard deviation; ASN, Asian; EUR, European; PRS, polygenic risk score.

490 **Table 2.** Association of polygenic risk scores and VTE risk in the test set individuals of European and African ancestry.

| Method | PRS tuning population | PRS | Number of SNPs | PRS testing population | | | |
|---|---|---|---|---|---|---|---|
| | | | | European | | African | |
| | | | | AUC (95% CI) | Odds ratio per SD (95% CI) | AUC (95% CI) | Odds ratio (95% CI) |
| (1) LDpred2 trained in EUR | - | LDpred2 $_{EUR}$ | 604,741 | 0.6 (0.59, 0.6) | 1.42 (1.39, 1.45) | 0.54 (0.53, 0.56) | 1.17 (1.11, 1.23) |
| (2) LDpred2 trained in AFR | - | LDpred2 $_{AFR}$ | 1,184,805 | 0.52 (0.52, 0.53) | 1.09 (1.06, 1.11) | 0.56 (0.54, 0.57) | 1.26 (1.2, 1.33) |
| Combine (1) + (2) | European | [a]LDpred2_combined $_{EUR}$ | 1,212,566 | 0.6 (0.59, 0.61) | 1.42 (1.39, 1.46) | 0.55 (0.53, 0.56) | 1.19 (1.13, 1.25) |
| Combine (1) + (2) | African | [a]LDpred2_combined $_{AFR}$ | 1,212,566 | 0.6 (0.59, 0.6) | 1.42 (1.39, 1.45) | 0.57 (0.55, 0.58) | 1.26 (1.2, 1.33) |
| (3) PRSCS trained in EUR | - | PRSCSX $_{EUR}$ | 591,788 | 0.61 (0.6, 0.61) | 1.47 (1.44, 1.51) | 0.56 (0.55, 0.57) | 1.26 (1.2, 1.33) |
| (4) PRSCS trained in AFR | - | PRSCSX $_{AFR}$ | 586,660 | 0.58 (0.58, 0.59) | 1.35 (1.32, 1.39) | 0.58 (0.57, 0.6) | 1.35 (1.28, 1.42) |
| Combine (3) + (4) | European | [a]PRSCSX_combined $_{EUR}$ | 598,977 | 0.61 (0.6, 0.62) | 1.48 (1.45, 1.52) | 0.57 (0.56, 0.59) | 1.3 (1.24, 1.37) |
| Combine (3) + (4) | African | [a]PRSCSX_combined $_{AFR}$ | 598,977 | 0.6 (0.6, 0.61) | 1.47 (1.44, 1.51) | 0.59 (0.57, 0.6) | 1.38 (1.3, 1.45) |

491 [a]Combined PRSs were generated using the formula $\alpha_0 + \alpha_1 PRS_1 + \alpha_2 PRS_2$ where $\alpha_0$, $\alpha_1$ and $\alpha_2$ are the weights obtained by fitting a logistic
492 regression model with VTE as outcome, PRS1 and PRS2 as explanatory variables using the validation data set. The weights for the considered
493 combination of PRSs can be found at https://github.com/yonhojee/VTE_PRS.