

# Genetically adjusted PSA levels for prostate cancer screening

Received: 14 March 2022

Accepted: 27 February 2023

Published online: 1 June 2023

 Check for updates

Linda Kachuri<sup>1,2,3</sup>, Thomas J. Hoffmann<sup>1,4</sup>, Yu Jiang<sup>2</sup>, Sonja I. Berndt<sup>5</sup>, John P. Shelley<sup>6</sup>, Kerry R. Schaffer<sup>7</sup>, Mitchell J. Machiela<sup>5</sup>, Neal D. Freedman<sup>5</sup>, Wen-Yi Huang<sup>5</sup>, Shengchao A. Li<sup>5</sup>, Ryder Easterlin<sup>8</sup>, Phyllis J. Goodman<sup>9</sup>, Cathee Till<sup>10</sup>, Ian Thompson<sup>11</sup>, Hans Lilja<sup>12,13</sup>, Stephen K. Van Den Eeden<sup>14</sup>, Stephen J. Chanock<sup>5</sup>, Christopher A. Haiman<sup>15,16</sup>, David V. Conti<sup>15,16</sup>, Robert J. Klein<sup>17</sup>, Jonathan D. Mosley<sup>6,18</sup>, Rebecca E. Graff<sup>1,20</sup> ✉ & John S. Witte<sup>1,2,3,19,20</sup> ✉

Prostate-specific antigen (PSA) screening for prostate cancer remains controversial because it increases overdiagnosis and overtreatment of clinically insignificant tumors. Accounting for genetic determinants of constitutive, non-cancer-related PSA variation has potential to improve screening utility. In this study, we discovered 128 genome-wide significant associations ( $P < 5 \times 10^{-8}$ ) in a multi-ancestry meta-analysis of 95,768 men and developed a PSA polygenic score (PGS<sub>PSA</sub>) that explains 9.61% of constitutive PSA variation. We found that, in men of European ancestry, using PGS-adjusted PSA would avoid up to 31% of negative prostate biopsies but also result in 12% fewer biopsies in patients with prostate cancer, mostly with Gleason score <7 tumors. Genetically adjusted PSA was more predictive of aggressive prostate cancer (odds ratio (OR) = 3.44,  $P = 6.2 \times 10^{-14}$ , area under the curve (AUC) = 0.755) than unadjusted PSA (OR = 3.31,  $P = 1.1 \times 10^{-12}$ , AUC = 0.738) in 106 cases and 23,667 controls. Compared to a prostate cancer PGS alone (AUC = 0.712), including genetically adjusted PSA improved detection of aggressive disease (AUC = 0.786,  $P = 7.2 \times 10^{-4}$ ). Our findings highlight the potential utility of incorporating PGS for personalized biomarkers in prostate cancer screening.

Prostate-specific antigen (PSA) is an enzyme produced by the prostate gland that degrades gel-forming seminal proteins to release motile sperm and is encoded by the *KLK3* (kallikrein 3) gene<sup>1–3</sup>. As prostate epithelial tissue becomes disrupted by a tumor, greater PSA concentrations are released into circulation<sup>2,3</sup>. PSA levels can also rise due to prostatic inflammation, infection, benign prostatic hyperplasia, older age and increased prostate volume<sup>3–5</sup>. Increased body mass index is associated with lower PSA levels, but the underlying mechanisms remain unclear<sup>6,7</sup>. Low PSA levels, thus, do not rule out prostate cancer, and PSA elevation is not sufficient for a conclusive diagnosis<sup>8</sup>. Although PSA testing reduces deaths from prostate cancer<sup>9</sup>, between

20% and 60% of cancers detected using PSA testing are estimated to be overdiagnoses<sup>10–12</sup>. In addition, the long-term risk of lethal prostate cancer remains low, especially in men with PSA below the age-specific median<sup>13,14</sup>. As a result, clinical guidelines in the United States and globally advise against population-level PSA screening and promote a shared decision-making model<sup>15,16</sup>.

One avenue for refining PSA screening is to account for variability in PSA due to genetic factors. PSA is highly heritable, with 40 independent loci identified in the largest previous genome-wide association study (GWAS)<sup>17,18</sup>. The goal of genetically correcting PSA levels is to increase the relative variation in PSA attributable to prostate cancer,

A full list of affiliations appears at the end of the paper. ✉ e-mail: [Rebecca.Graff@ucsf.edu](mailto:Rebecca.Graff@ucsf.edu); [jswitte@stanford.edu](mailto:jswitte@stanford.edu)

thereby improving their predictive value for disease detection. The first study to genetically correct PSA using just four variants reclassified 3% of participants to warranting biopsy and 3% to avoiding biopsy<sup>19</sup>. Incorporating additional genetic predictors has the potential to personalize PSA testing, reduce overdiagnosis-related morbidity and improve detection of lethal disease. To maximize the utility of this approach, it is critical to distinguish genetic variants that influence constitutive PSA levels from those affecting prostate tumor development. PSA and prostate cancer share many genetic loci<sup>17,19–22</sup>, but the extent to which this overlap reflects screening bias remains unclear, as GWASs of prostate cancer may capture signals for disease susceptibility and incidental detection due to benign PSA elevation.

Our study explores the genetic architecture of PSA levels in men without prostate cancer, with a view toward assessing whether genetic adjustment of PSA improves clinical decision-making related to prostate cancer diagnosis. It also provides a novel framework for the clinical translation of polygenic scores (PGSs) for non-causal cancer biomarkers.

## Results

The study design of the Precision PSA study is illustrated in Fig. 1. Using data from five studies (Methods), we conducted genome-wide analyses of PSA levels  $\leq 10$  ng ml<sup>-1</sup> in cis-gender men never diagnosed with prostate cancer. GWAS results were meta-analyzed within ancestry groups and then combined across populations for a total sample size of 95,768 individuals.

### Genetic architecture of PSA variation

The heritability ( $h^2$ ) of PSA levels was investigated using several methods to assess sensitivity to underlying modeling assumptions (Methods). Across 26,491 men of European ancestry in the UK Biobank (UKB) with linked clinical records, the median PSA value was 2.35 ng ml<sup>-1</sup> (Supplementary Fig. 1). Using individual-level data for variants with minor allele frequency (MAF)  $\geq 0.01$  and imputation INFO  $> 0.80$ , PSA heritability was  $h^2 = 0.41$  (95% confidence interval (CI): 0.36–0.46) based on GCTA<sup>23</sup> and  $h^2 = 0.30$  (95% CI: 0.26–0.33) based on LDAK<sup>24</sup> (Supplementary Table 1 and Extended Data Fig. 1). Applying LDAK to GWAS summary statistics generated from the same individuals produced similar estimates ( $h^2 = 0.35$ , 95% CI: 0.28–0.43), whereas other methods<sup>25,26</sup> were biased downward. In the European ancestry GWAS meta-analysis ( $n_{\text{EUR}} = 85,824$ ), LDAK estimated  $h^2 = 0.30$  (95% CI: 0.29–0.31). Sample sizes for other ancestries were too small for reliable heritability estimates.

The multi-ancestry meta-analysis of 95,768 men from five studies identified 128 independent index variants ( $P < 5.0 \times 10^{-8}$ , linkage disequilibrium (LD)  $r^2 < 0.01$  within  $\pm 10$ -Mb windows) across 90 chromosomal cytoband regions (Fig. 2). The strongest associations were in known PSA loci<sup>17,19,21,22</sup>, such as *KLK3* (rs17632542,  $P = 3.2 \times 10^{-638}$ ), 10q26.12 (rs10886902,  $P = 8.2 \times 10^{-118}$ ), *MSMB* (rs10993994,  $P = 7.3 \times 10^{-87}$ ), *NKX3-1* (rs1160267,  $P = 6.3 \times 10^{-83}$ ), *CLPTMIL* (rs401681,  $P = 7.0 \times 10^{-54}$ ) and *HNF1B* (rs10908278,  $P = 2.1 \times 10^{-46}$ ). Eighty-two index variants were independent of previously detected associations in the Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort<sup>17</sup>; they mapped to 56 cytobands where PSA signals have not previously been reported. Associations initially detected in the UKB (Extended Data Fig. 1b) strengthened in the meta-analysis: *TEX11* in Xq13.1 (rs62608084,  $P = 1.7 \times 10^{-24}$ ); *THADA* in 2p21 (rs11899863,  $P = 1.7 \times 10^{-13}$ ); *OTX1* in 2p15 (rs58235267,  $P = 4.9 \times 10^{-13}$ ); *SALL3* in 18q23 (rs71279357,  $P = 1.8 \times 10^{-12}$ ); and *ST6GAL1* in 3q27.3 (rs12629450,  $P = 2.6 \times 10^{-10}$ ). Additional novel findings included *CDK5RAP1* (rs291671,  $P = 1.2 \times 10^{-18}$ ), *LDAH* (rs10193919,  $P = 1.5 \times 10^{-15}$ ), *ABCC4* (rs61965887,  $P = 3.7 \times 10^{-14}$ ), *INKA2* (rs2076591,  $P = 2.6 \times 10^{-13}$ ), *SUDS3* (rs1045542,  $P = 1.2 \times 10^{-13}$ ), *FAF1* (rs12569177,  $P = 3.2 \times 10^{-13}$ ), *JARID2* (rs926309,  $P = 1.6 \times 10^{-12}$ ), *GPC3* (rs4829762,  $P = 5.9 \times 10^{-12}$ ), *EDA* (rs2520386,  $P = 4.2 \times 10^{-11}$ ) and *ODF3* (rs7103852,  $P = 1.2 \times 10^{-9}$ ) (Supplementary Tables 2 and 3).

Of the 128 index variants, 96 reached genome-wide significance in the European ancestry meta-analysis, as did three in the East Asian ancestry meta-analysis (*KLK3*: rs2735837 and rs374546878; *MSMB*: rs10993994;  $n_{\text{EAS}} = 3,337$ ), two in the Hispanic/Latino meta-analysis (*KLK3*: rs17632542 and rs2735837;  $n_{\text{HIS/LAT}} = 3,098$ ) and one in the African ancestry meta-analysis (*FGFR2*: rs10749415;  $n_{\text{AFR}} = 3,509$ ) (Supplementary Table 4). Effect sizes from the European ancestry GWAS were modestly correlated with estimates from other ancestries (Spearman's  $\rho_{\text{HIS/LAT}} = 0.48$ ,  $P = 1.1 \times 10^{-8}$ ;  $\rho_{\text{AFR}} = 0.27$ ,  $P = 2.0 \times 10^{-3}$ ;  $\rho_{\text{EAS}} = 0.16$ ,  $P = 0.068$ ) (Supplementary Fig. 2). However, cross-population comparisons of correlations should be interpreted with caution as they are confounded by higher sampling error in groups with smaller sample sizes.

There was heterogeneity (Cochran's  $Q$   $P_Q < 0.05$ ) across ancestry-specific fixed-effects meta-analyses for 12 of 128 index variants, four of which had effects in different directions: rs58235267 (*OTX1*), rs1054713 (*KLK1*), rs10250340 (*EIF4HPI*) and rs7020681 (*SLC35D2*) (Supplementary Table 5). An alternative meta-analysis approach, MR-MEGA<sup>27</sup>, which partitions effect size heterogeneity into components correlated with ancestry and residual variation, identified one additional signal in 5q15 (rs291812,  $P_{\text{MR-MEGA}} = 1.0 \times 10^{-8}$ ) that was driven by the East Asian ancestry results ( $P_{\text{EAS}} = 1.2 \times 10^{-6}$ ) (Supplementary Table 6).

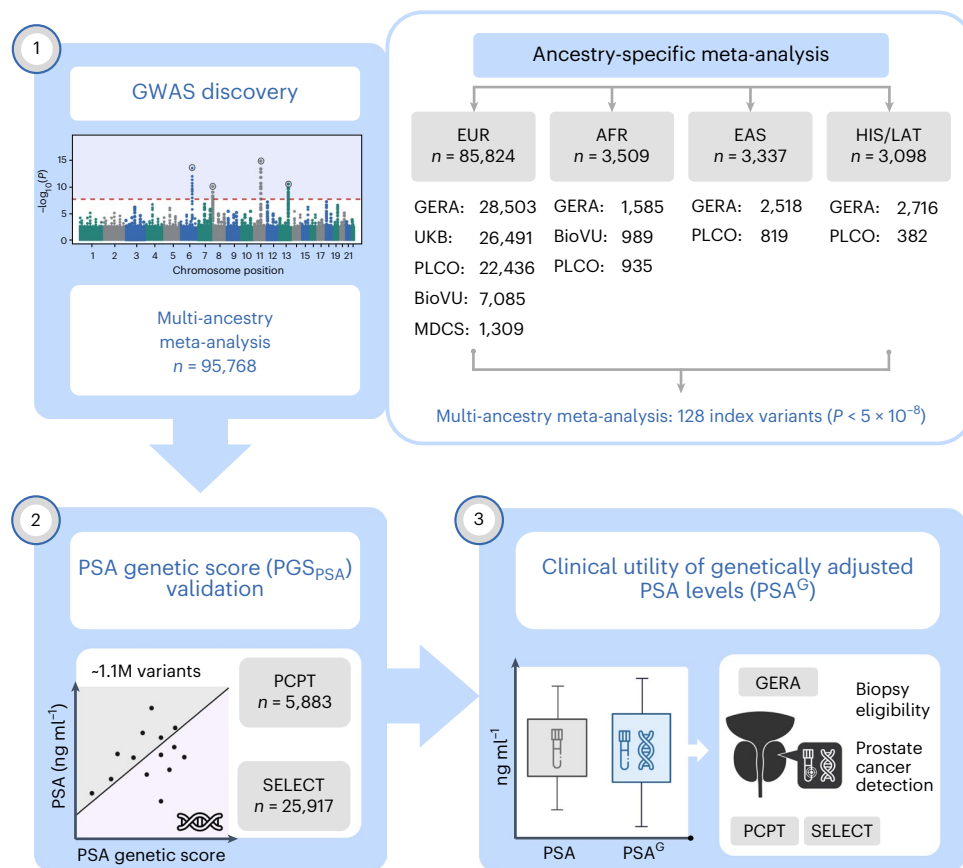
Predicted functional consequences of the 128 index variants were explored using CADD<sup>28</sup>. Scores  $> 13$  (corresponding to the 5% most deleterious substitutions genome-wide) were observed for 16 of the 128 index variants detected in the original fixed effects meta-analysis, including ten new signals: rs10193919 (*LDAH*); rs7732515 in 5q14.3; rs11899863 (*THADA*); rs58235267 (*OTX1*); rs926309 (*JARID2*); rs4829762 (*GPC3*) and rs13268, a missense variant in *FBN1*; rs78378222 in *TP53* and rs3760230 in *SMG6*; and rs712329 in *SLC25A21* (Supplementary Table 7). Sixty-one variants had significant (false discovery rate (FDR)  $< 0.05$ ) effects on gene expression, including 15 prostate tissue expression quantitative trait loci (eQTLs) for 17 eGenes, 55 blood eQTLs for 185 eGenes and nine eQTLs with effects in both tissues. Notable eGenes included *RUVBL1*, a chromatin-remodeling factor that modulates pro-inflammatory NF- $\kappa$ B signaling and transcription of Myc and  $\beta$ -catenin<sup>29</sup>; *ODF3*, which maintains elastic structures in the sperm tail<sup>30</sup>; and *LDAH*, which promotes cholesterol mobilization in macrophages<sup>31</sup>. Several PSA-associated variants were eQTLs for genes involved in immune response (*IFITM2*, *IFITM3* and *HS1BP3*).

### Impact of PSA-related selection bias on prostate cancer GWAS

Because prostate cancer detection often hinges on PSA elevation, genetic factors resulting in higher constitutive PSA levels may appear to increase prostate cancer risk because of more frequent screening. Of the 128 lead PSA variants, 52 (41%) were associated with prostate cancer at the Bonferroni-corrected threshold ( $P < 0.05/128$ ) in the PRACTICAL consortium's European ancestry GWAS<sup>32</sup> (Supplementary Table 8). Using the method by Dudbridge et al.<sup>33</sup>, we investigated whether index event bias could partly explain these shared signals<sup>33,34</sup> (Methods, Fig. 3 and Supplementary Table 9). Applying the estimated bias correction factor ( $b = 1.144$ ) decreased the number of variants associated with prostate cancer from 52 to 34 (Extended Data Fig. 2). When we corrected 209 European ancestry prostate cancer risk variants ( $P < 5.0 \times 10^{-8}$ , LD  $r^2 < 0.01$ ) for screening bias, 93 (45%) remained genome-wide significant. Notably, rs76765083 (*KLK3*) remained genome-wide significant but reversed direction. Sensitivity analyses using SlopeHunter<sup>35</sup> resulted in 150 (72%) variants with  $P < 5 \times 10^{-8}$  (Supplementary Table 10).

### Development and validation of PGS<sub>PSA</sub>

We considered two approaches for constructing a PGS for PSA: clumping genome-wide significant associations from the multi-ancestry meta-analysis (PGS<sub>128</sub>) and a genome-wide score generated using the Bayesian PRS-CSx algorithm (PGS<sub>CSx</sub>) (ref.<sup>36</sup>) (Methods). Each score was validated in the Prostate Cancer Prevention Trial (PCPT)



**Fig. 1 | Overview of the Precision PSA study design.** Genome-wide association analyses were conducted in men without prostate cancer and meta-analyzed within each population: European ancestry (EUR), African ancestry (AFR), East Asian ancestry (EAS) and Hispanic/Latino ancestry (HIS/LAT). Ancestry-stratified results were used to develop a genome-wide  $PGS_{PSA}$  comprised of approximately

1.1 million variants and were also combined into a multi-ancestry meta-analysis of 95,768 men.  $PGS_{PSA}$  was validated in the PCPT and the SELECT and was used to compute  $PSA^G$  values. We examined how using  $PSA^G$  values affects eligibility for prostate biopsy and evaluated associations with incident prostate cancer.

and the Selenium and Vitamin E Cancer Prevention Trial (SELECT), which were excluded from the discovery GWAS. Most of the men in both cohorts were of European ancestry, although SELECT offered larger sample sizes for other ancestry groups (Extended Data Fig. 3).  $PGS_{CSx}$  was ultimately selected, as it was more predictive of baseline PSA than  $PGS_{128}$  in multi-ancestry analyses and most ancestry subgroups (Supplementary Table 11).

In the PCPT,  $PGS_{CSx}$  accounted for 8.13% of variation in baseline PSA levels ( $\beta$  per s.d. increase = 0.186,  $P = 3.3 \times 10^{-112}$ ) in the pooled multi-ancestry sample of 5,883 men (Fig. 4a–c and Supplementary Table 11).  $PGS_{CSx}$  was associated with PSA across age groups, although effects attenuated in participants aged  $\geq 70$  years (Extended Data Fig. 4).  $PGS_{CSx}$  was validated in 5,725 participants of European ancestry ( $EUR \geq 0.80$ ) ( $PGS_{CSx}$ :  $\beta = 0.194$ ,  $P = 1.7 \times 10^{-115}$ ), but neither  $PGS_{128}$  nor  $PGS_{CSx}$  reached nominal significance in the admixed European and African ancestry ( $0.20 < AFR/EUR < 0.80$ ,  $n = 103$ ) or East Asian ancestry ( $EAS \geq 0.80$ ,  $n = 55$ ) populations.

In the SELECT,  $PGS_{CSx}$  was associated with baseline PSA levels in the pooled sample of 25,917 men ( $\beta = 0.258$ ,  $P = 1.3 \times 10^{-619}$ ) and among men of European ancestry ( $n = 22,253$ ,  $\beta_{PGS} = 0.283$ ,  $P = 5.5 \times 10^{-610}$ ), accounting for 9.61% to 10.94% of variation, respectively (Fig. 4b–d and Supplementary Table 11).  $PGS_{CSx}$  also validated in the East Asian ( $n = 257$ ,  $\beta = 0.258$ ,  $P = 5.9 \times 10^{-7}$ ) and admixed EAS/EUR ( $n = 321$ ,  $\beta = 0.315$ ,  $P = 5.2 \times 10^{-12}$ ) ancestry groups. In men with admixed AFR/EUR ancestry ( $n = 1,763$ ),  $PGS_{CSx}$  explained 4.22% of PSA variation ( $\beta = 0.157$ ,  $P = 4.8 \times 10^{-19}$ ).  $PGS_{128}$  was more predictive than  $PGS_{CSx}$  ( $\beta = 0.163$ ,  $P = 8.2 \times 10^{-11}$  versus  $\beta = 0.098$ ,  $P = 8.0 \times 10^{-6}$ ) in men of African ancestry

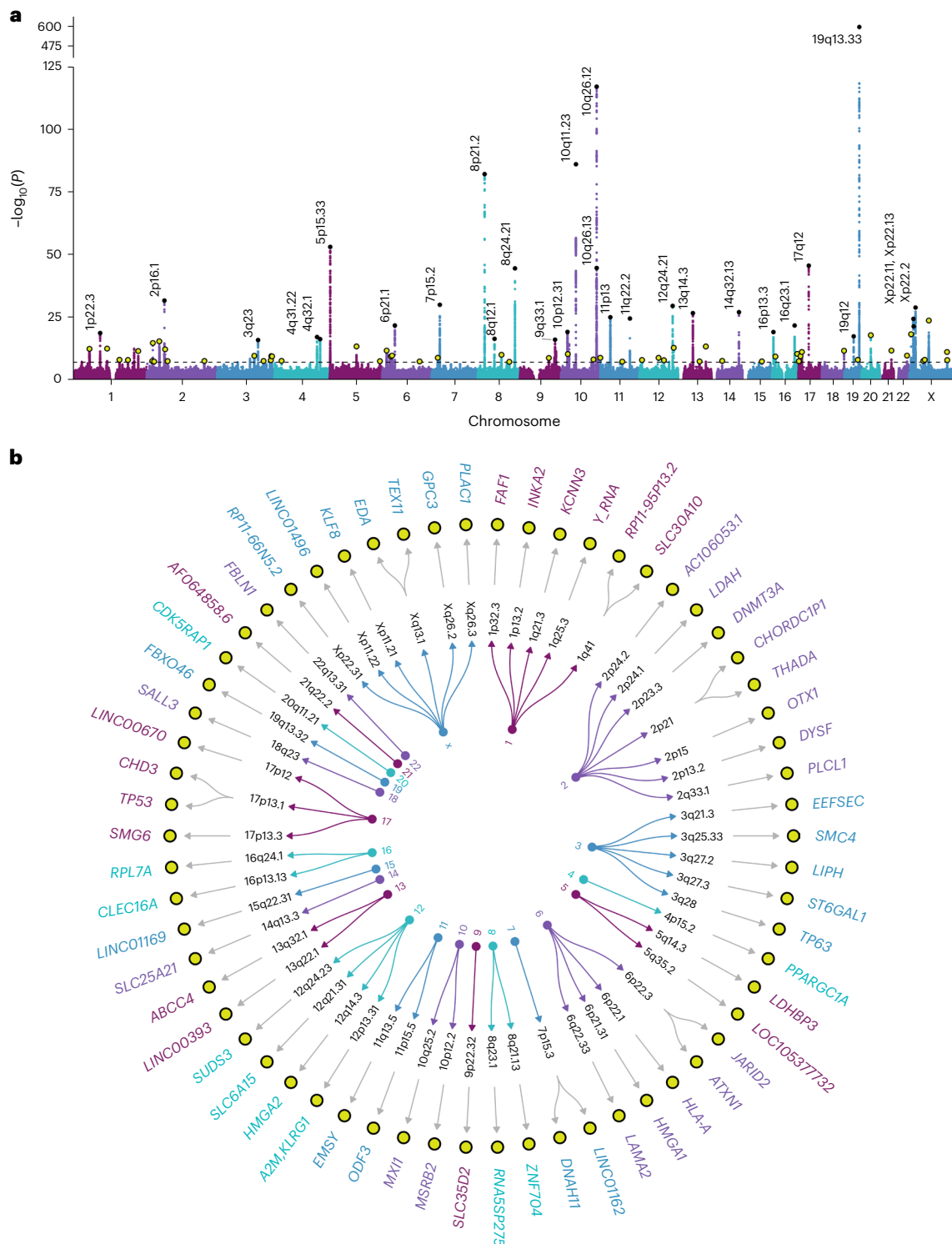
( $AFR \geq 0.80$ ,  $n = 1,173$ ) and the pooled AFR and admixed ( $0.20 < EUR/AFR < 0.80$ ) group ( $n = 2,936$ ).

We also examined associations with temporal trends in PSA: velocity, calculated using  $\log(PSA)$  values at two timepoints, and doubling time in months (Methods and Supplementary Table 12). In men with a PSA increase (SELECT pooled sample:  $n = 14,908$ ),  $PGS_{CSx}$  was associated with less rapid velocity ( $PGS_{CSx}$ :  $\beta = -4.06 \times 10^{-4}$ ,  $P = 3.7 \times 10^{-5}$ ) and longer doubling time ( $PGS_{CSx}$ :  $\beta = 10.41$ ,  $P = 1.9 \times 10^{-8}$ ). In men with a PSA decrease between the first and last timepoint (SELECT pooled sample:  $n = 6,970$ ),  $PGS_{CSx}$  was only suggestively associated with slowing PSA decline ( $\beta = 5.02 \times 10^{-4}$ ,  $P = 0.068$ ). The same pattern was observed in the PCPT, with higher  $PGS_{CSx}$  values conferring less rapid changes in PSA.

$PGS_{CSx}$ , referred to as  $PGS_{PSA}$  from here onward, was used to genetically adjust baseline or earliest pre-randomization PSA values ( $PSA^G$ ) for each individual, relative to the population mean (Methods and equations 1 and 2).  $PSA^G$  and unadjusted PSA were strongly correlated in the PCPT (Pearson's  $r = 0.841$ ,  $0.833$ – $0.848$ ) and the SELECT ( $r = 0.854$ ,  $0.851$ – $0.857$ ). The number of participants with  $PSA^G > 4$  ng ml<sup>-1</sup>, a commonly used threshold for diagnostic testing, increased from 0 to 24 in the PCPT and from 5 to 413 in the SELECT (Fig. 4e,f), reflecting the preferential trial selection of men with low PSA<sup>8,37</sup>.

### Impact of PSA-related bias on PGS associations

In men of European ancestry in the UKB excluded from the PSA GWAS, there was a strong positive relationship between the 269-variant prostate cancer PGS ( $PGS_{269}$ )<sup>32</sup> and  $PGS_{PSA}$  in cases ( $n = 11,568$ ,  $\beta = 0.190$ ,



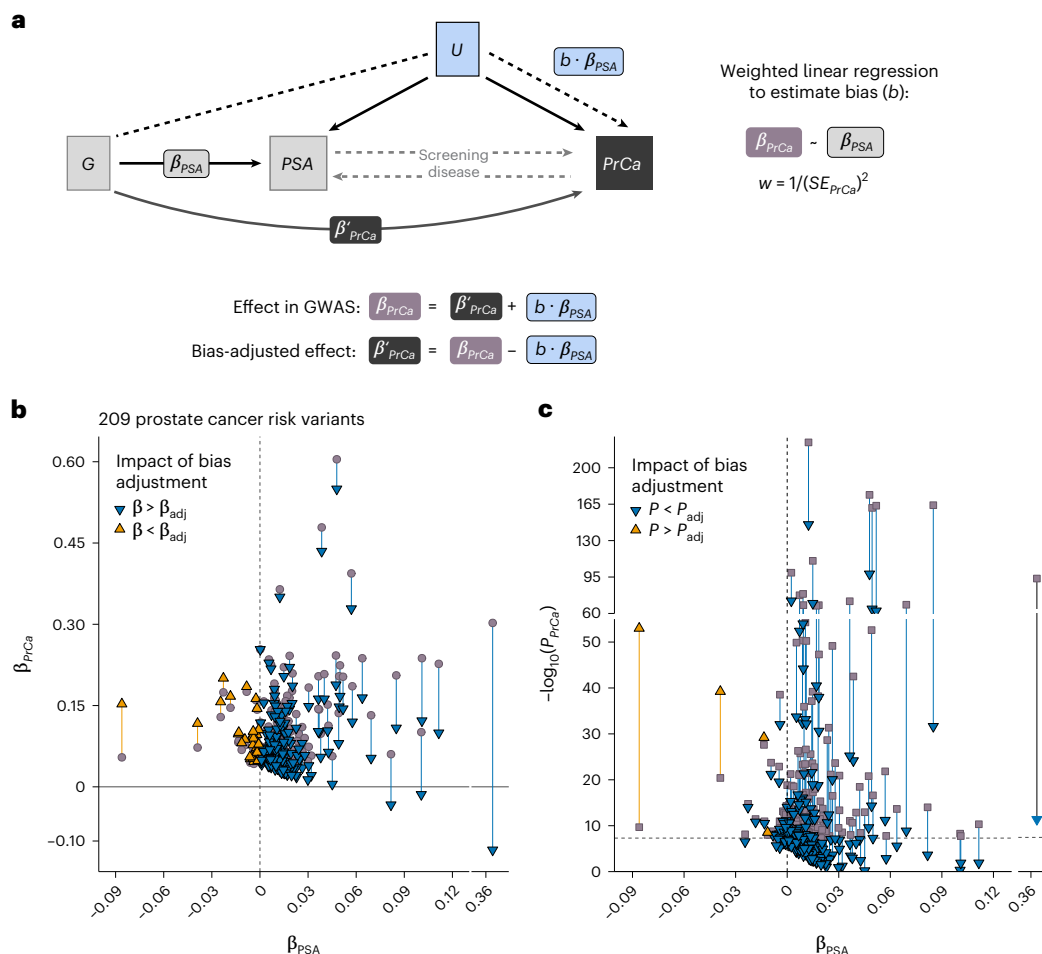
**Fig. 2 | Multi-ancestry GWAS of PSA levels. a**, Manhattan plot depicting the results of the GWAS meta-analysis of PSA levels in 95,768 men without prostate cancer. The genome-wide significance threshold of  $P < 5 \times 10^{-8}$  is indicated by the dotted black line. Index variants within known PSA-associated loci are annotated with the corresponding cytoband. Novel findings are highlighted in yellow.

**b**, Circular dendrogram shows the nearest gene(s) for novel PSA-associated variants. Genome-wide significant ( $P < 5 \times 10^{-8}$ ) index variants were selected using LD-based clumping ( $LD R^2 < 0.01$  within  $\pm 10$ -Mb windows). All GWAS  $P$  values are two-sided and derived from a fixed-effects inverse-variance-weighted meta-analysis using METAL.

$P = 2.3 \times 10^{-96}$ ) and controls ( $n = 152,884$ ,  $\beta = 0.236$ ,  $P < 10^{-700}$ ) (Extended Data Fig. 5 and Supplementary Table 13). Re-fitting  $PGS_{269}$  using weights corrected for index event bias ( $PGS_{269}^{adj}$ ) substantially attenuated associations in cases ( $\beta_{adj} = 0.029$ ,  $P = 2.7 \times 10^{-3}$ ) and controls ( $\beta_{adj} = 0.052$ ,  $P = 2.2 \times 10^{-89}$ ).

To further characterize the impact of this bias, we examined  $PGS_{269}$  associations with prostate cancer status in 3,673 cases and 2,363 biopsy-confirmed, European ancestry controls from GERA.  $PGS_{269}^{adj}$  had a larger magnitude of association with prostate cancer (OR for top decile = 3.63, 95% CI: 3.01–4.37) than  $PGS_{269}$  (odds ratio (OR) = 2.71, 95%





**Fig. 3 | Influence of PSA-related index event bias on prostate cancer**

**GWAS. a**, Conceptual diagram depicts how selection on PSA levels induces an association between genetic variant  $G$  and  $U$ , a composite confounder that captures polygenic and non-genetic factors. This selection induces an association with prostate cancer (PrCa) via path  $G \rightarrow U \rightarrow PrCa$ , in addition to the direct  $G \rightarrow PrCa$  effects. Bi-directional dotted lines show that PSA is not only a disease biomarker but also influences screening behavior and the likelihood of prostate cancer detection. **b, c**, The impact of bias correction is shown for 209 prostate cancer risk variants. Independent risk variants were selected from the PRACTICAL GWAS meta-analysis (85,554 cases and 91,972 controls of European

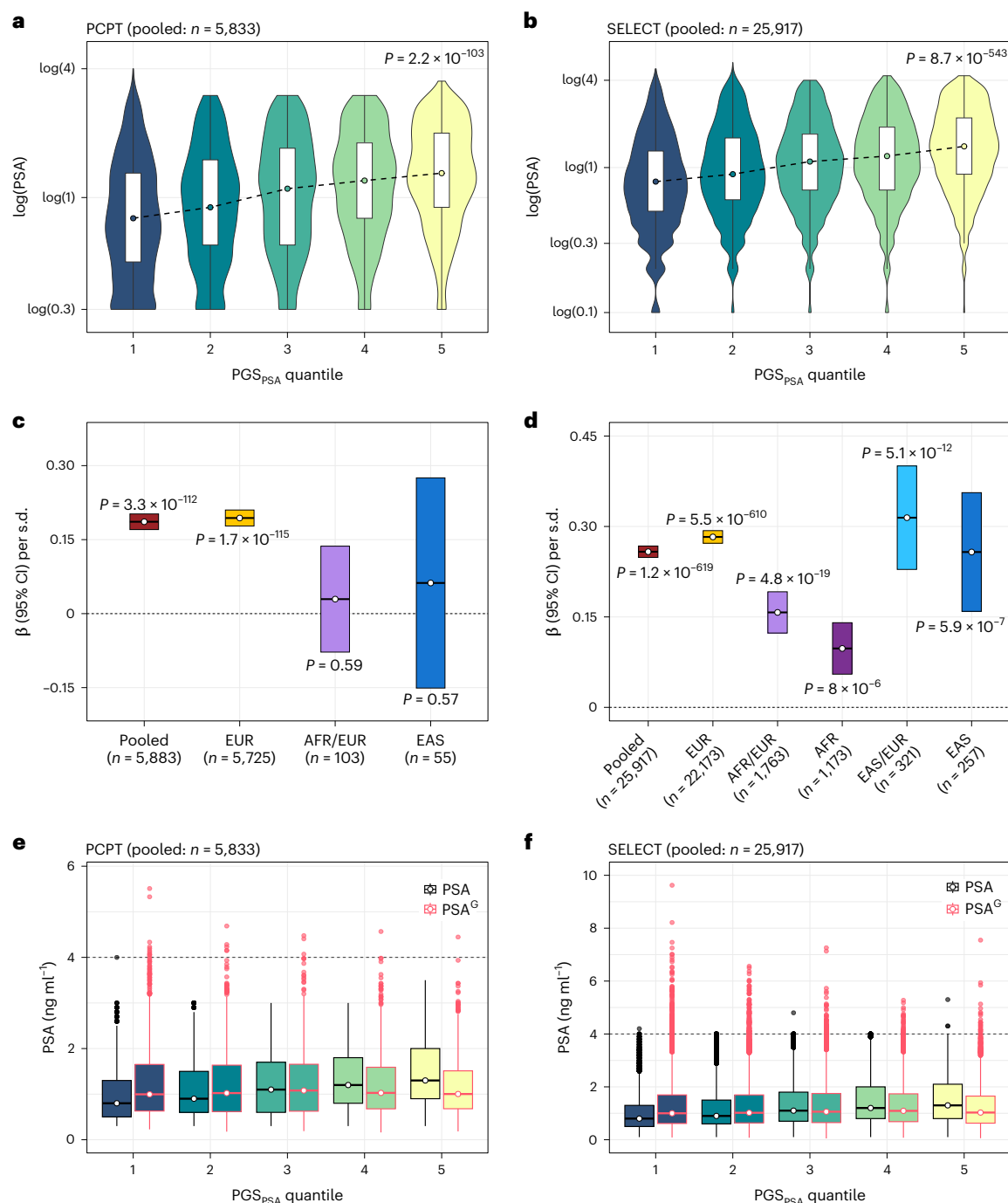
ancestry) by Conti et al.<sup>32</sup> using LD clumping ( $LD R^2 < 0.01$ ,  $P < 5 \times 10^{-8}$ ). For each variant, associations with PSA ( $\beta_{PSA}$ ) are based on an inverse-variance-weighted fixed-effects meta-analysis in men of European ancestry ( $n = 85,824$ ). **b**, GWAS effect sizes for prostate cancer ( $\beta_{PrCa}$ ) are aligned to the risk-increasing allele. Bias-adjusted effect sizes ( $\beta_{adj}$ ) are denoted by triangles. **c**, Two-sided GWAS  $P$  values for prostate cancer ( $P_{PrCa}$ ) were derived from an inverse-variance-weighted fixed-effects meta-analysis. Two-sided bias-adjusted  $P$  values ( $P_{adj}$ ), denoted by triangles, were calculated from a chi-squared test statistic based on  $\beta_{adj}$  and corresponding standard errors. Genome-wide significance threshold ( $P < 5 \times 10^{-8}$ ) is indicated by the horizontal dotted line.

CI: 2.28–3.21) and higher area under the curve (AUC: 0.685 versus 0.677,  $P = 3.91 \times 10^{-3}$ ) (Supplementary Table 14). The impact of bias correction was most pronounced for Gleason  $\geq 7$  tumors ( $PGS_{269}^{adj}$  AUC = 0.692 versus  $PGS_{269}$  AUC = 0.678,  $P = 1.91 \times 10^{-3}$ ), although these AUC estimates are inflated due to overlap with the GWAS used to develop  $PGS_{269}$  (ref. 32). In case-only analyses,  $PGS_{PSA}$  and  $PGS_{269}$  were inversely associated with Gleason score, illustrating how screening bias decreases the likelihood of identifying high-grade disease (Supplementary Table 15). Compared to Gleason  $\leq 6$  tumors, an s.d. increase in  $PGS_{PSA}$  was inversely associated with Gleason 7 disease (OR = 0.79, 95% CI: 0.76–0.83) and Gleason  $\geq 8$  disease (OR = 0.71, 95% CI: 0.64–0.81). Patients in the top decile of  $PGS_{269}$  were approximately 30% less likely to have Gleason  $\geq 8$  tumors (OR = 0.72, 95% CI: 0.54–0.96) than Gleason  $\leq 6$  tumors, but this association was attenuated after bias correction ( $PGS_{269}^{adj}$ : OR = 0.94, 95% CI: 0.75–1.17).

### Impact of genetic adjustment of PSA on biopsy eligibility

Among GERA participants who underwent prostate biopsy, we examined how adjustment using  $PGS_{PSA}$  reclassified individuals for biopsy

recommendation at age-specific thresholds used by Kaiser Permanente: 40–49 years old = 2.5 ng ml<sup>-1</sup>; 50–59 years old = 3.5 ng ml<sup>-1</sup>; 60–69 years old = 4.5 ng ml<sup>-1</sup>; and 70–79 years old = 6.5 ng ml<sup>-1</sup> (Methods). For men of European ancestry, mean PSA levels in men with a negative biopsy ( $n = 2,363$ , 7.2 ng ml<sup>-1</sup>) were higher than in men without prostate cancer who did not have a biopsy ( $n = 24,811$ , 1.5 ng ml<sup>-1</sup>) (Supplementary Table 16). Relative to all controls, where standardized  $PGS_{PSA} = 0$ , biopsied men were enriched for PSA-increasing alleles (cases:  $PGS_{PSA} = 0.278$ ; controls:  $PGS_{PSA} = 0.934$ ). After genetic adjustment, 31.7% of biopsy-negative men were reclassified below the PSA level for recommending biopsy, and 2.5% became biopsy eligible, resulting in a net reclassification of 29.3% (27.5% to 31.21%) (Fig. 5a). Among 3,673 cases, PSA<sup>G</sup> values below the biopsy referral threshold were more prevalent than upward adjustment, resulting in a net reclassification of -8.6% (-9.48% to -7.67%) (Fig. 5a). Of the patients who became ineligible, most had Gleason < 7 tumors ( $n = 300$ , 72%; Supplementary Table 16). In men of African ancestry, there were few changes in biopsy eligibility among patients ( $n = 392$ ), with 3.1% reclassified upward and 4.6% downward (Fig. 5b and Supplementary Table 16). Of 108 biopsy-negative



**Fig. 4 | Validation of the  $\text{PGS}_{\text{PSA}}$  in two cancer prevention trials.**

**a–f**, Performance of  $\text{PGS}_{\text{PSA}}$  was evaluated in the PCPT and the SELECT. **a,b**, Violin plots show the distribution of baseline  $\log(\text{PSA})$  within quantiles of  $\text{PGS}_{\text{PSA}}$ , comprising 1,058,173 and 1,071,278 variants in the PCPT (**a**) and the SELECT (**b**). Box plots extend from the 25th to the 75th percentiles, with a trend line connecting the median value within each age stratum. Two-sided  $P$  values were derived from linear regression models for the effect of a quantile increase in  $\text{PGS}_{\text{PSA}}$  on  $\log(\text{PSA})$ . **c,d**, Crossbar plots show the effect estimates ( $\beta$ ) and corresponding 95% CIs per s.d. increase in the standardized  $\text{PGS}_{\text{PSA}}$  on baseline

$\log(\text{PSA})$  in the PCPT (**c**) and the SELECT (**d**). Ancestry-stratified and pooled multi-ancestry estimates are presented. Two-sided  $P$  values based on linear regression models are annotated. **e,f**, Comparison of distributions for PSA and  $\text{PSA}^{\text{G}}$ , with the horizontal line at  $4 \text{ ng ml}^{-1}$ , a commonly used threshold for further diagnostic testing. Box plots show the median value, with lower and upper hinges corresponding to the 25th and 75th percentiles or first and third quartiles. Whiskers extend as a multiple of the interquartile range ( $\text{IQR} \times 1.5$ ). Outlying values beyond the end of the whiskers are plotted individually.

controls, 75 (69.4%) were reclassified below the referral threshold based on  $\text{PSA}^{\text{G}}$ , reflecting high enrichment for predisposition to PSA elevation ( $\text{PGS}_{\text{PSA}} = 1.710$ ). The overall net reclassification was positive, suggesting that  $\text{PSA}^{\text{G}}$  has some clinical utility in both populations.

### PSA genetic adjustment improves prostate cancer detection

The utility of  $\text{PSA}^{\text{G}}$ , alone and in combination with  $\text{PGS}_{269}$ , was first assessed in the PCPT, where end-of-study biopsies were performed in all participants, effectively eliminating potential misclassification

of prostate cancer status. Among 335 cases and 5,548 controls,  $\text{PGS}_{\text{PSA}}$  was not associated with prostate cancer incidence (pooled: OR per s.d. = 1.01,  $P = 0.83$ ), confirming that it captures genetic determinants of non-cancer PSA variation. The magnitude of association for genetically adjusted baseline  $\text{PSA}^G$  with prostate cancer (OR per unit increase in  $\log(\text{PSA ng ml}^{-1}) = 1.90$ , 95% CI: 1.56–2.31) was slightly larger than for PSA (OR = 1.88, 95% CI: 1.55–2.29) in the European ancestry group (Supplementary Table 17). The magnitude of association with prostate cancer was larger for  $\text{PGS}_{269}^{\text{adj}}$  (pooled and European: OR per s.d. = 1.57, 95% CI: 1.40–1.76) than for  $\text{PGS}_{269}$  without bias correction (pooled: OR = 1.52, 95% CI: 1.36–1.70; European: OR = 1.53, 95% CI: 1.36–1.72) (Supplementary Table 17). The model with  $\text{PGS}_{269}^{\text{adj}}$  and  $\text{PSA}^G$  achieved the best classification in the pooled (AUC = 0.686) and European ancestry (AUC = 0.688) populations and outperformed  $\text{PGS}_{269}^{\text{adj}}$  alone (pooled: AUC = 0.656,  $P_{\text{AUC}} = 7.5 \times 10^{-4}$ ; European: AUC = 0.658,  $P_{\text{AUC}} = 1.4 \times 10^{-3}$ ).

The benefit of genetically adjusting PSA was most evident for detection of aggressive prostate cancer, defined as Gleason  $\geq 7$ , PSA  $\geq 10 \text{ ng ml}^{-1}$ , T3–T4 stage and/or distant or nodal metastases. In the PCPT,  $\text{PSA}^G$  conferred an approximately threefold risk increase (pooled: OR = 2.87, 95% CI: 1.98–4.65, AUC = 0.706; European: OR = 2.99, 95% CI: 1.95–4.59, AUC = 0.711) compared to  $\text{PGS}_{269}^{\text{adj}}$  (pooled: OR = 1.55, 95% CI: 1.23–1.95, AUC = 0.651; European: OR = 1.55, 95% CI: 1.22–1.96, AUC = 0.657) (Fig. 6a and Supplementary Table 18). The model with  $\text{PSA}^G$  and  $\text{PGS}_{269}^{\text{adj}}$  achieved AUC = 0.726 (European: AUC = 0.734) for aggressive tumors but had lower discrimination for non-aggressive disease (pooled and European: AUC = 0.681) (Supplementary Table 19). Among patients with prostate cancer,  $\text{PSA}^G$  (pooled: OR = 2.06, 95% CI: 1.23–3.45) and baseline PSA (pooled: OR = 1.81, 95% CI: 1.12–3.10) were associated with higher likelihood of aggressive compared to non-aggressive tumors, whereas  $\text{PGS}_{269}$  (pooled: OR = 0.91,  $P = 0.54$ ) and  $\text{PGS}_{269}^{\text{adj}}$  (OR = 0.97,  $P = 0.85$ ) were not (Supplementary Table 20).

In the SELECT, associations with risk of prostate cancer overall (Supplementary Table 21), aggressive disease (Fig. 6b and Supplementary Table 22) and non-aggressive disease (Supplementary Table 23) in the pooled and European ancestry analyses were similar to the PCPT. In men of East Asian ancestry, associations for  $\text{PSA}^G$  (OR = 2.15, 95% CI: 0.82–5.62) were attenuated compared to PSA (OR = 2.60, 95% CI: 1.03–6.54). This was also observed in men of African ancestry, although the effect size for  $\text{PSA}^G$  derived using  $\text{PGS}_{128}$  (OR = 3.37, 95% CI: 2.38–4.78) was larger than for  $\text{PSA}^G$  based on  $\text{PGS}_{\text{CSX}}$  (OR = 2.68, 95% CI: 1.94–3.69), consistent with the larger proportion of variation in PSA explained by  $\text{PGS}_{128}$  than  $\text{PGS}_{\text{CSX}}$  in this population. Models for prostate cancer including  $\text{PSA}^G$  were calibrated in the pooled and European ancestry individuals, whereas, in the African ancestry subgroup,  $\text{PSA}^G$  inaccurately estimated risk in upper deciles (Supplementary Figs. 3–6).

The largest improvement in discrimination from  $\text{PSA}^G$  (OR = 3.81, 95% CI: 2.62–5.54, AUC = 0.777) relative to PSA (OR = 3.40, 95% CI: 2.34–4.93, AUC = 0.742,  $P_{\text{AUC}} = 0.026$ ) and to  $\text{PGS}_{269}$  (OR = 1.76, 95% CI: 1.41–2.21, AUC = 0.726,  $P_{\text{AUC}} = 0.057$ ) was for aggressive tumors in men of European ancestry (106 cases, 23,667 controls). In the pooled African ancestry population (18 cases, 2,733 controls),  $\text{PSA}^G$  based on  $\text{PGS}_{128}$  (OR = 2.96, 95% CI: 1.43–6.12), but not  $\text{PGS}_{\text{CSX}}$  (OR = 2.48, 95% CI: 1.24–4.97), was more predictive than unadjusted PSA (OR = 2.82, 95% CI: 1.33–5.99) (Supplementary Table 22). The best model for aggressive disease included  $\text{PSA}^G$  and  $\text{PGS}_{269}^{\text{adj}}$  for pooled (AUC = 0.788, 95% CI: 0.744–0.831) and European ancestry (AUC = 0.804, 95% CI: 0.757–0.851) populations, but, for African ancestry individuals, unadjusted PSA and  $\text{PGS}_{269}$  without bias correction achieved the highest AUC of 0.828 (95% CI: 0.739–0.916).  $\text{PSA}^G$  was better calibrated than PSA in pooled and European ancestry groups but not in African ancestry participants (Supplementary Figs. 7 and 8).

## Discussion

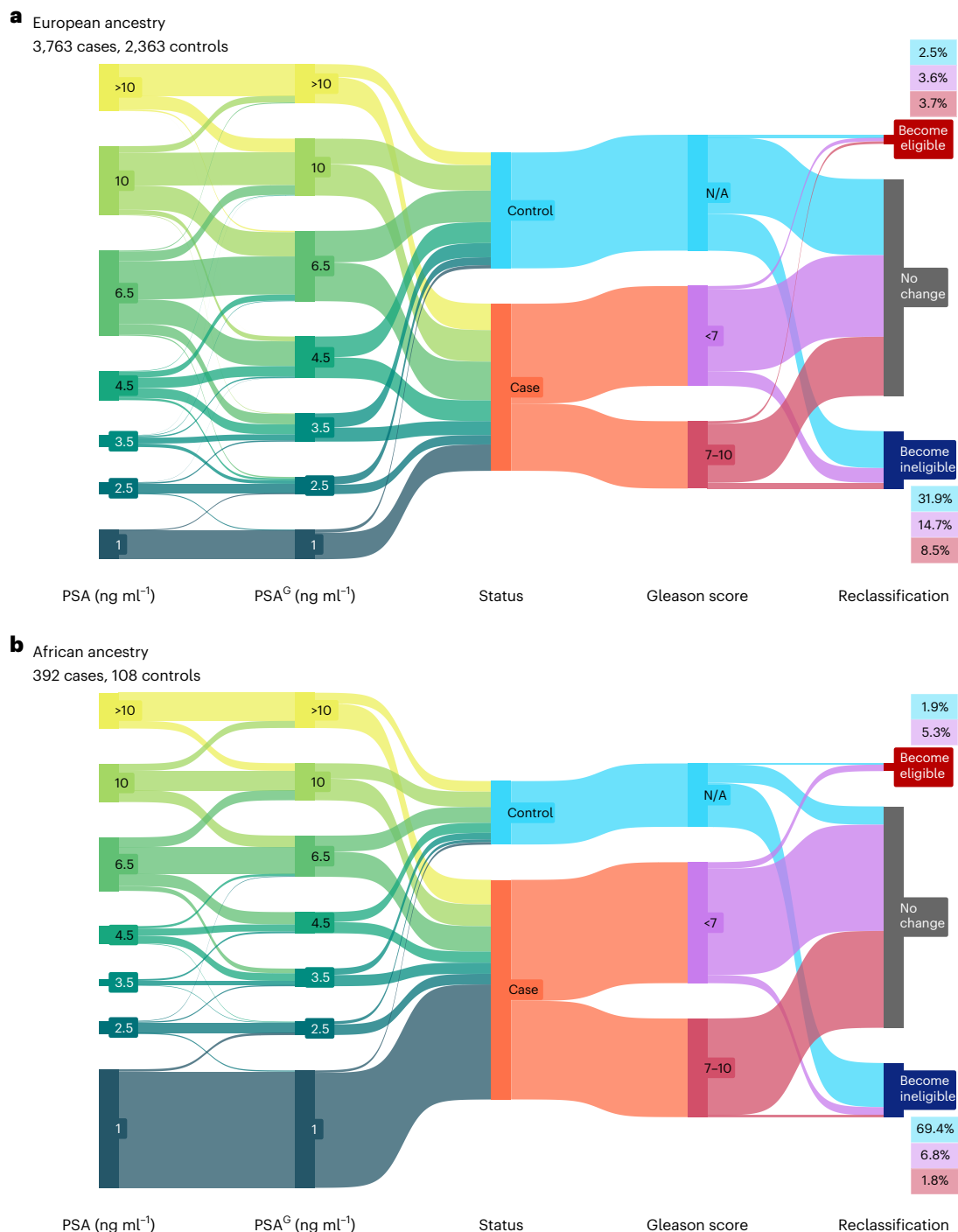
Serum PSA is the most widely used biomarker for prostate cancer detection, although concerns with specificity and, to a lesser degree,

sensitivity have limited adoption of PSA testing for population-level screening. Leveraging PGS to personalize diagnostic biomarkers, such as PSA, provides a new avenue for translating GWAS discoveries into clinical practice. This concept, termed ‘de-Mendelization’, is essentially Mendelian randomization in reverse—subtracting the genetically predicted component of trait variance instead of using it to estimate causal effects. De-Mendelization of non-causal predictive biomarkers can maximize disease-related signal and improve disease detection<sup>38,39</sup>. Although previous work on PSA genetics<sup>19</sup> and other biomarkers<sup>38,40</sup> has alluded to the potential of genetic adjustment to produce clinically meaningful shifts in the PSA distribution, the value of this approach for reducing overdiagnosis and detecting aggressive disease has not been previously shown.

Risk-stratified, personalized screening for prostate cancer will require parallel efforts to elucidate the genetic architecture of prostate cancer susceptibility and PSA variation in individuals without disease. Our GWAS advances these efforts by discovering 82 novel PSA-associated variants. The strongest novel signals map to genes involved in reproductive processes, potentially reflecting non-cancer function of PSA in liquefying seminal fluid. *TEX11* on Xq13.1, for example, is preferentially expressed in male germ cells and early spermatocytes. *TEX11* mutations cause meiotic arrest and azoospermia, and this gene regulates homologous chromosome synapsis and double-strand DNA break repair<sup>41</sup>. *ODF3* encodes a component of sperm flagella fibers and has been linked to regulation of platelet count and volume<sup>42</sup>. Other novel loci contained genes involved in embryonic development, epigenetic regulation and chromatin organization, including *DNMT3A*, *OTX1*, *CHD3*, *JARID2*, *HMG1A*, *HMG2* and *SUSD3*. *DNMT3A* is a methyltransferase that regulates imprinting and X-chromosome inactivation and has been studied extensively in the context of height<sup>43</sup>, clonal hematopoiesis and hematologic cancers<sup>44</sup>. *CHD3* is involved in chromatin remodeling during development and suppresses herpes simplex virus infection<sup>45</sup>. Multiple PSA-associated variants were in genes related to infection and immunity, including *HLA-A*; *ST6GAL1*, involved in IgG N-glycosylation<sup>46</sup>; *KLRG1*, which regulates natural killer (NK) cell function and IFN- $\gamma$  production<sup>47</sup>; and *FUT2*, which affects ABO precursor H antigen presentation and confers susceptibility to viral and bacterial infections<sup>48</sup>.

Although our GWAS was restricted to men without prostate cancer, several cancer susceptibility genes were among the PSA-associated loci, including a pan-cancer risk variant in *TP53* (rs78378222) (ref.<sup>49</sup>) and signals in *TP63*, *GPC3* and *THADA*. Although we cannot rule out undiagnosed prostate cancer in our participants, its prevalence is unlikely to be high enough to produce appreciable bias. Pervasive pleiotropy and omnigenic architecture<sup>50</sup> may explain the diverse functions of PSA loci implicated in inflammation, epigenetic regulation and growth factor signaling. Even established tumor suppressor genes, such as *TP53*, *GPC3* and *THADA*, have pleiotropic effects on obesity via dysregulation of cell growth and metabolism<sup>51–53</sup>. Furthermore, distinct p63 isoforms regulate epithelial and craniofacial development as well as apoptosis of male germ cells and spermatogenesis<sup>54,55</sup>. Mutations in *GPC3* cause Simson–Golabi–Behmel syndrome, which is characterized by visceral and skeletal abnormalities and excess risk of embryonic tumors<sup>56</sup>.

Distinguishing variants that influence prostate cancer detection via PSA screening from genetic signals for prostate carcinogenesis has implications for deciphering biological mechanisms and developing risk prediction models. Prostate cancer detection depends on PSA testing, whereas PSA screening is influenced by genetic factors affecting constitutive PSA levels. The bias arising from this complex relationship may be substantial. Our findings suggest that bias-corrected effect sizes more accurately capture the contribution of GWAS-identified variants to prostate cancer risk, without conflating it with detection. Correction for PSA-related bias and subsequent improvement in  $\text{PGS}_{269}$  performance for detecting aggressive disease is an extension of de-Mendelization. Adjusting risk allele weights may be a more effective strategy than filtering out variants based on associations with PSA.



**Fig. 5 | Genetically adjusted PSA influences biopsy eligibility.** **a–b**, Flow diagrams illustrate changes in PSA values after genetic adjustment for participants in the GERA cohort and subsequent reclassification at PSA thresholds used to recommend prostate biopsy. Genetic adjustment was applied to the last pre-biopsy PSA value to obtain PSA<sup>G</sup>. Analyses were performed

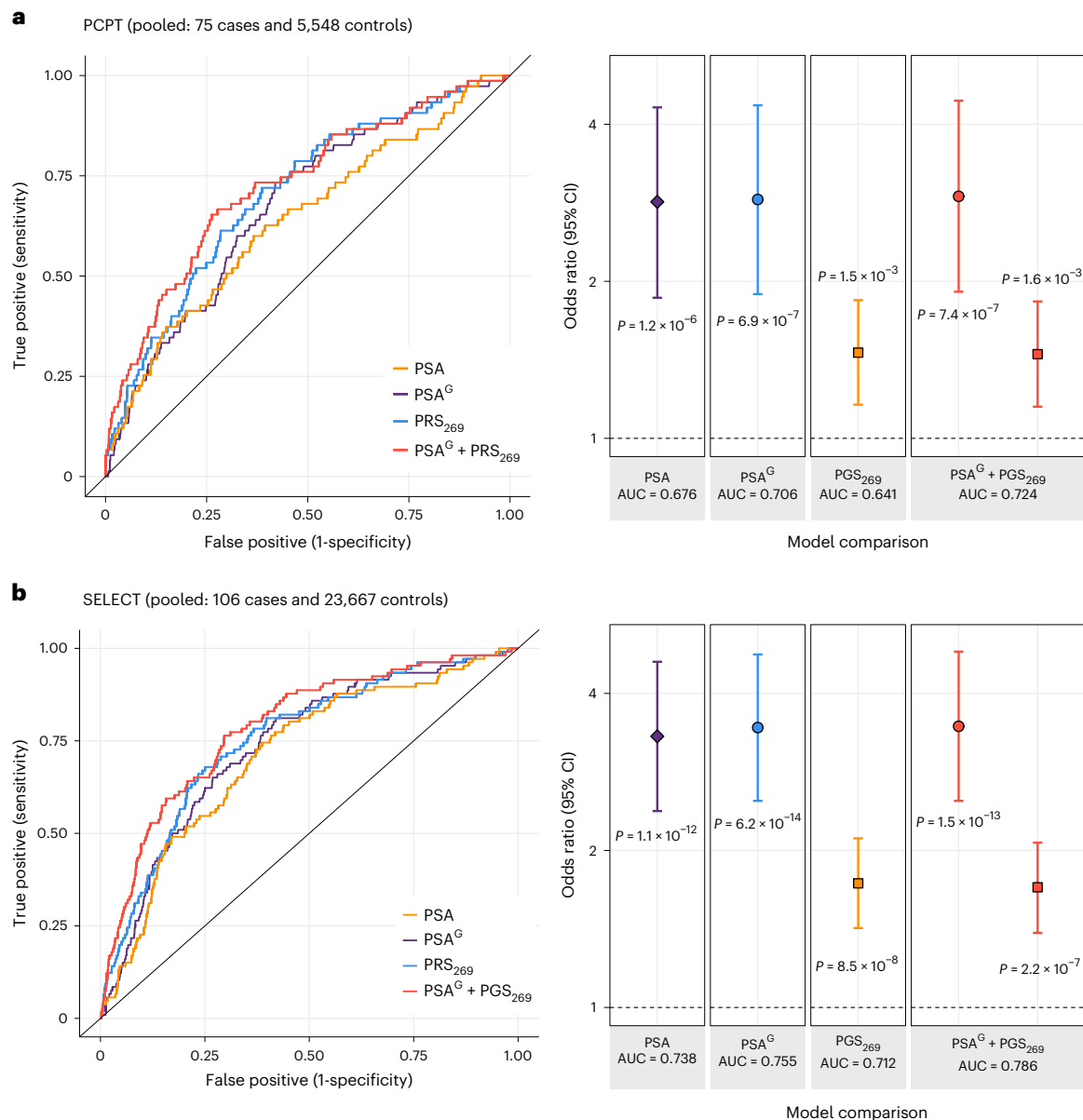
separately in men of European (**a**) and African (**b**) ancestry. Size of the nodes and flows are proportional to the number of individuals in each category. Patients with prostate cancer (cases) were stratified by Gleason score categories, where Gleason <7 represents potentially indolent disease. Gleason score is not applicable to men with a negative prostate biopsy (controls).

Generally, the improvements in PSA<sup>G</sup> and PGS<sub>269</sub> are proportional to the extent of their de-noising of signals for PSA elevation unrelated to prostate cancer. The impact of bias correction was most pronounced in populations selected for high PSA, such as men who underwent prostate biopsy in GERA, but it was also observed in the PCPT and the SELECT, which enrolled men with low PSA.

Our investigation of index event bias has several limitations. The Dudbridge method assumes that direct genetic effects on PSA

levels and prostate cancer susceptibility are uncorrelated, and violations of this assumption over-attribute shared genetic signals to selection bias<sup>33</sup>. Although SlopeHunter relaxes this assumption<sup>35</sup>, analyses of PGS<sub>269</sub> suggest that it under-corrects selection bias. SlopeHunter relies on clustering to distinguish PSA-specific from pleiotropic variants<sup>35</sup>, with small or poorly separated clusters resulting in unstable bias estimates. Disentangling genetic associations between PSA and prostate cancer with greater certainty will





**Fig. 6 | Genetic associations with aggressive prostate cancer. a–b,** Comparison of models for aggressive disease, defined as Gleason score  $\geq 7$ , PSA  $\geq 10$  ng ml<sup>-1</sup>, T3–T4 stage and/or distant or nodal metastases in the PCPT (**a**) and the SELECT (**b**). The pooled study population includes all ancestry groups. Logistic regression models were adjusted for baseline age, randomization arm, the top ten population-specific genetic ancestry principal components and proportions

of African and East Asian genetic ancestry. ORs and 95% CIs were estimated per 1-unit increase in log(PSA ng ml<sup>-1</sup>) and log(PSA<sup>G</sup> ng ml<sup>-1</sup>) and per s.d. increase in the prostate cancer genetic risk score (PGS<sub>269</sub>) from Conti et al.<sup>32</sup>, which was standardized to achieve s.d. equal to 1. All *P* values are two-sided. AUC is based on the full model with all covariates.

require experiments such as CRISPR screens and massively parallel reporter assays.

Another limitation is that the reported magnitude of biopsy reclassification may be specific to GERA and Kaiser Permanente clinical guidelines and biased because GERA controls comprised 30% of the PSA discovery GWAS. Because it was unlikely for men with low PSA to be biopsied, and most patients with prostate cancer already had PSA values at or above the biopsy referral cutoff, there were limited opportunities to increase biopsy eligibility in this population. Despite these limitations, our findings indicate that genetically adjusted PSA may reduce overdiagnosis and overtreatment, albeit accompanied by some undesirable loss of sensitivity. Although reclassifying cases to not receive biopsy is concerning, most such reclassifications occurred among patients with non-aggressive disease, a group susceptible to overdiagnosis<sup>37</sup>.

Our PGS-based approach updates the first application of PSA genetic correction by Gudmundsson et al.<sup>19</sup> while retaining straightforward calculation of the genetic correction factor. Increasing the specificity of an established, clinically useful biomarker is efficient and would have low adoption barriers. However, analytic choices, such as selecting an optimal PGS algorithm and reference population for obtaining mean PGS<sub>PSA</sub>, are not trivial. The choice of reference population affects the magnitude of correction and clinical decisions based on absolute PSA values. Furthermore, any new biomarker would require validation in real-world settings to identify populations who would benefit most and characterize barriers to implementation, such as physician familiarity with PGS and patient education about genetic testing. Genetically adjusted PSA should also be evaluated in conjunction with other procedures used for prostate cancer detection, such as

targeted magnetic resonance imaging, and explored as a criterion for refining selection of participants into screening trials.

Our study highlights the importance and challenge of developing a PGS that adequately performs across the spectrum of ancestry. Compared to PGS<sub>128</sub>, PGS<sub>CSX</sub> did not improve performance in men of African ancestry. This may reflect the ‘meta’ estimation procedure, which does not require a separate dataset for hyperparameter tuning but is less accurate<sup>36</sup>. GWAS efforts in larger and more diverse cohorts are underway and will expand the catalog of PSA-associated variants and increase their utility. Genetic adjustment using a PGS<sub>PSA</sub> that does not explain a sufficiently high proportion of trait variation risks decreasing the accuracy of PSA screening.

Future research should assess whether genetically adjusted PSA levels improve prediction of prostate cancer mortality and investigate PSA-related biomarkers, such as the ratio of free to total PSA and pro-PSA (a precursor PSA isoform), which may have higher specificity for prostate cancer detection<sup>58,59</sup>. Although PGS<sub>PSA</sub> was associated with PSA doubling time and velocity, these metrics assess change between two timepoints and may not capture PSA trajectories that are meaningful for disease detection<sup>60</sup>. Clinical guidelines for PSA kinetics are also lacking in the context of prostate cancer screening. Regardless, we think that genetic adjustment may improve the accuracy of any heritable PSA biomarker and may be a valuable addition to multi-omic biomarkers.

In summary, by detecting genetic variants associated with non-prostate cancer PSA variation, we developed a PGS<sub>PSA</sub> that captures the contribution of common genetic variants to a man’s inherent PSA level. We showed that a straightforward calculation of genetically adjusted, personalized PSA levels using PGS<sub>PSA</sub> provides clinically meaningful improvements in prostate cancer diagnostic characteristics. Moreover, genetic determinants of PSA provide an avenue for mitigating selection bias due to PSA screening in prostate cancer GWASs and improving disease prediction. These results illustrate a proof of concept for incorporating genetic factors into PSA screening for prostate cancer and expanding this approach to other diagnostic biomarkers.

## Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41591-023-02277-9>.

## References

- Lilja, H. A kallikrein-like serine protease in prostatic fluid cleaves the predominant seminal vesicle protein. *J. Clin. Invest.* **76**, 1899–1903 (1985).
- Balk, S. P., Ko, Y. J. & Bubley, G. J. Biology of prostate-specific antigen. *J. Clin. Oncol.* **21**, 383–391 (2003).
- Lilja, H., Ulmert, D. & Vickers, A. J. Prostate-specific antigen and prostate cancer: prediction, detection and monitoring. *Nat. Rev. Cancer* **8**, 268–278 (2008).
- Pinsky, P. F. et al. Prostate volume and prostate-specific antigen levels in men enrolled in a large screening trial. *Urology* **68**, 352–356 (2006).
- Lee, S. E. et al. Relationship of prostate-specific antigen and prostate volume in Korean men with biopsy-proven benign prostatic hyperplasia. *Urology* **71**, 395–398 (2008).
- Grubb, R. L. 3rd et al. Serum prostate-specific antigen hemodilution among obese men undergoing screening in the Prostate, Lung, Colorectal, and Ovarian Cancer Screening Trial. *Cancer Epidemiol. Biomark. Prev.* **18**, 748–751 (2009).
- Harrison, S. et al. Systematic review and meta-analysis of the associations between body mass index, prostate cancer, advanced prostate cancer, and prostate-specific antigen. *Cancer Causes Control* **31**, 431–449 (2020).
- Thompson, I. M. et al. Assessing prostate cancer risk: results from the Prostate Cancer Prevention Trial. *J. Natl Cancer Inst.* **98**, 529–534 (2006).
- Schroder, F. H. et al. Screening and prostate-cancer mortality in a randomized European study. *N. Engl. J. Med.* **360**, 1320–1328 (2009).
- Telesca, D., Etzioni, R. & Gulati, R. Estimating lead time and overdiagnosis associated with PSA screening from prostate cancer incidence trends. *Biometrics* **64**, 10–19 (2008).
- Welch, H. G. & Black, W. C. Overdiagnosis in cancer. *J. Natl Cancer Inst.* **102**, 605–613 (2010).
- Vickers, A. J. et al. Empirical estimates of prostate cancer overdiagnosis by age and prostate-specific antigen. *BMC Med.* **12**, 26 (2014).
- Vickers, A. J. et al. Strategy for detection of prostate cancer based on relation between prostate specific antigen at age 40–55 and long term risk of metastasis: case–control study. *BMJ* **346**, f2023 (2013).
- Kovac, E. et al. Association of baseline prostate-specific antigen level with long-term diagnosis of clinically significant prostate cancer among patients aged 55 to 60 years: a secondary analysis of a cohort in the Prostate, Lung, Colorectal, and Ovarian (PLCO) Cancer Screening Trial. *JAMA Netw. Open* **3**, e1919284 (2020).
- Tikkinen, K. A. O. et al. Prostate cancer screening with prostate-specific antigen (PSA) test: a clinical practice guideline. *BMJ* **362**, k3581 (2018).
- US Preventive Services Task Force et al. Screening for prostate cancer: US Preventive Services Task Force recommendation statement. *JAMA* **319**, 1901–1913 (2018).
- Hoffmann, T. J. et al. Genome-wide association study of prostate-specific antigen levels identifies novel loci independent of prostate cancer. *Nat. Commun.* **8**, 14248 (2017).
- Bansal, A. et al. Heritability of prostate-specific antigen and relationship with zonal prostate volumes in aging twins. *J. Clin. Endocrinol. Metab.* **85**, 1272–1276 (2000).
- Gudmundsson, J. et al. Genetic correction of PSA values using sequence variants associated with PSA levels. *Sci. Transl. Med.* **2**, 62ra92 (2010).
- Benafif, S., Kote-Jarai, Z., Eeles, R. A. & Consortium, P. A review of prostate cancer genome-wide association studies (GWAS). *Cancer Epidemiol. Biomark. Prev.* **27**, 845–857 (2018).
- Wiklund, F. et al. Association of reported prostate cancer risk alleles with PSA levels among men without a diagnosis of prostate cancer. *Prostate* **69**, 419–427 (2009).
- Kim, S., Shin, C. & Jee, S. H. Genetic variants at 1q32.1, 10q11.2 and 19q13.41 are associated with prostate-specific antigen for prostate cancer screening in two Korean population-based cohort studies. *Gene* **556**, 199–205 (2015).
- Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- Speed, D., Holmes, J. & Balding, D. J. Evaluating and improving heritability models using summary statistics. *Nat. Genet.* **52**, 458–462 (2020).
- Bulik-Sullivan, B. K. et al. LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
- Ning, Z., Pawitan, Y. & Shen, X. High-definition likelihood inference of genetic correlations across human complex traits. *Nat. Genet.* **52**, 859–864 (2020).
- Magi, R. et al. Trans-ethnic meta-regression of genome-wide association studies accounting for ancestry increases power for discovery and improves fine-mapping resolution. *Hum. Mol. Genet.* **26**, 3639–3650 (2017).

28. Rentzsch, P., Schubach, M., Shendure, J. & Kircher, M. CADD-Splice—improving genome-wide variant effect prediction using deep learning-derived splice scores. *Genome Med.* **13**, 31 (2021).
29. Mao, Y. Q. & Houry, W. A. The role of pontin and reptin in cellular physiology and cancer etiology. *Front. Mol. Biosci.* **4**, 58 (2017).
30. Egydio de Carvalho, C. et al. Molecular cloning and characterization of a complementary DNA encoding sperm tail protein SHIPPO 1. *Biol. Reprod.* **66**, 785–795 (2002).
31. Currall, B. B. et al. Loss of LDAH associated with prostate cancer and hearing loss. *Hum. Mol. Genet.* **27**, 4194–4203 (2018).
32. Conti, D. V. et al. Trans-ancestry genome-wide association meta-analysis of prostate cancer identifies new susceptibility loci and informs genetic risk prediction. *Nat. Genet.* **53**, 65–75 (2021).
33. Dudbridge, F. et al. Adjustment for index event bias in genome-wide association studies of subsequent events. *Nat. Commun.* **10**, 1561 (2019).
34. Paternoster, L., Tilling, K. & Davey Smith, G. Genetic epidemiology and Mendelian randomization for informing disease therapeutics: conceptual and methodological challenges. *PLoS Genet.* **13**, e1006944 (2017).
35. Mahmoud, O., Dudbridge, F., Davey Smith, G., Munafo, M. & Tilling, K. A robust method for collider bias correction in conditional genome-wide association studies. *Nat. Commun.* **13**, 619 (2022).
36. Ruan, Y. et al. Improving polygenic prediction in ancestrally diverse populations. *Nat. Genet.* **54**, 573–580 (2022).
37. Lippman, S. M. et al. Effect of selenium and vitamin E on risk of prostate cancer and other cancers: the Selenium and Vitamin E Cancer Prevention Trial (SELECT). *JAMA* **301**, 39–51 (2009).
38. Kjaergaard, A. D., Bojesen, S. E., Nordestgaard, B. G., Johansen, J. S. & Smith, G. D. Biomarker de-Mendelization: principles, potentials and limitations of a strategy to improve biomarker prediction by reducing the component of variance explained by genotype. Preprint at *bioRxiv* <https://doi.org/10.1101/428276> (2018).
39. Holmes, M. V. & Davey Smith, G. Can Mendelian randomization shift into reverse gear? *Clin. Chem.* **65**, 363–366 (2019).
40. Enroth, S., Johansson, A., Enroth, S. B. & Gyllenstein, U. Strong effects of genetic and lifestyle factors on biomarker variation and use of personalized cutoffs. *Nat. Commun.* **5**, 4684 (2014).
41. Yatsenko, A. N. et al. X-linked *TEX11* mutations, meiotic arrest, and azoospermia in infertile men. *N. Engl. J. Med.* **372**, 2097–2107 (2015).
42. Astle, W. J. et al. The allelic landscape of human blood cell trait variation and links to common complex disease. *Cell* **167**, 1415–1429 (2016).
43. Marouli, E. et al. Rare and low-frequency coding variants alter human adult height. *Nature* **542**, 186–190 (2017).
44. Bick, A. G. et al. Inherited causes of clonal haematopoiesis in 97,691 whole genomes. *Nature* **586**, 763–768 (2020).
45. Arbuckle, J. H. & Kristie, T. M. Epigenetic repression of herpes simplex virus infection by the nucleosome remodeler CHD3. *mBio* **5**, e01027-13 (2014).
46. Shen, X. et al. Multivariate discovery and replication of five novel loci associated with immunoglobulin G N-glycosylation. *Nat. Commun.* **8**, 447 (2017).
47. Wang, J. M. et al. KLRG1 negatively regulates natural killer cell functions through the Akt pathway in individuals with chronic hepatitis C virus infection. *J. Virol.* **87**, 11626–11636 (2013).
48. Kachuri, L. et al. The landscape of host genetic factors involved in immune response to common viral infections. *Genome Med.* **12**, 93 (2020).
49. Rashkin, S. R. et al. Pan-cancer study detects genetic risk variants and shared genetic basis in two large cohorts. *Nat. Commun.* **11**, 4423 (2020).
50. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An expanded view of complex traits: from polygenic to omnigenic. *Cell* **169**, 1177–1186 (2017).
51. Di Giovannantonio, M. et al. Heritable genetic variants in key cancer genes link cancer risk with anthropometric traits. *J. Med. Genet.* **58**, 392–399 (2021).
52. Filmus, J. & Capurro, M. The role of glypican-3 in the regulation of body size and cancer. *Cell Cycle* **7**, 2787–2790 (2008).
53. Moraru, A. et al. THADA regulates the organismal balance between energy storage and heat production. *Dev. Cell* **41**, 72–81 e76 (2017).
54. Vanbokhoven, H., Melino, G., Candi, E. & Declercq, W. p63, a story of mice and men. *J. Invest. Dermatol.* **131**, 1196–1207 (2011).
55. Wang, H. et al. Transcriptional regulation of P63 on the apoptosis of male germ cells and three stages of spermatogenesis in mice. *Cell Death Dis.* **9**, 76 (2018).
56. Neri, G., Gurrieri, F., Zanni, G. & Lin, A. Clinical and molecular aspects of the Simpson–Golabi–Behmel syndrome. *Am. J. Med. Genet.* **79**, 279–283 (1998).
57. Gulati, R., Inoue, L. Y., Gore, J. L., Katcher, J. & Etzioni, R. Individualized estimates of overdiagnosis in screen-detected prostate cancer. *J. Natl Cancer Inst.* **106**, djt367 (2014).
58. Catalona, W. J. et al. Use of the percentage of free prostate-specific antigen to enhance differentiation of prostate cancer from benign prostatic disease: a prospective multicenter clinical trial. *JAMA* **279**, 1542–1547 (1998).
59. Loeb, S. et al. The prostate health index selectively identifies clinically significant prostate cancer. *J. Urol.* **193**, 1163–1169 (2015).
60. Vickers, A. J. & Brewster, S. F. PSA velocity and doubling time in diagnosis and prognosis of prostate cancer. *Br. J. Med Surg. Urol.* **5**, 162–168 (2012).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023, corrected publication 2025

<sup>1</sup>Department of Epidemiology & Biostatistics, University of California, San Francisco, San Francisco, CA, USA. <sup>2</sup>Department of Epidemiology & Population Health, Stanford University School of Medicine, Stanford, CA, USA. <sup>3</sup>Stanford Cancer Institute, Stanford University School of Medicine, Stanford, CA, USA. <sup>4</sup>Institute of Human Genetics, University of California, San Francisco, San Francisco, CA, USA. <sup>5</sup>Division of Cancer Epidemiology and Genetics, National Cancer Institute, Rockville, MD, USA. <sup>6</sup>Department of Biomedical Informatics, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>7</sup>Vanderbilt-Ingram Cancer Center, Nashville, TN, USA. <sup>8</sup>Biological and Medical Informatics, University of California, San Francisco, San Francisco, CA, USA.

<sup>9</sup>Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>10</sup>SWOG Statistics and Data Management Center, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. <sup>11</sup>CHRISTUS Santa Rosa Medical Center Hospital, San Antonio, TX, USA. <sup>12</sup>Departments of Laboratory Medicine, Surgery and Medicine, Memorial Sloan Kettering Cancer Center, New York, NY, USA. <sup>13</sup>Department of Translational Medicine, Lund University, Skåne University Hospital, Malmö, Sweden. <sup>14</sup>Division of Research, Kaiser Permanente Northern California, Oakland, CA, USA. <sup>15</sup>Center for Genetic Epidemiology, Department of Population and Preventive Health Sciences, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>16</sup>Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. <sup>17</sup>Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA. <sup>18</sup>Department of Internal Medicine, Vanderbilt University Medical Center, Nashville, TN, USA. <sup>19</sup>Departments of Biomedical Data Science and Genetics, Stanford University, Stanford, CA, USA. <sup>20</sup>These authors jointly supervised this work: Rebecca E. Graff, John S. Witte. ✉e-mail: [Rebecca.Graff@ucsf.edu](mailto:Rebecca.Graff@ucsf.edu); [jswitte@stanford.edu](mailto:jswitte@stanford.edu)



## Methods

Informed consent was obtained from all study participants. The UKB received ethics approval from the Research Ethics Committee (reference: 11/NW/0382) in accordance with the UKB Ethics and Governance Framework. The research was conducted with approved access to UKB data under application number 14105. We used previously published PSA GWAS results from the GERA cohort by Hoffmann et al.<sup>17</sup>. The original study was approved by the Kaiser Permanente Northern California institutional review board and the University of California, San Francisco Human Research Protection Program Committee on Human Research. The Prostate, Lung, Colorectal and Ovarian (PLCO) Cancer Screening Trial was approved by the institutional review board at each participating center and the National Cancer Institute. The informed consent document signed by PLCO study participants allows use of these data by investigators for discovery and hypothesis generation in the investigation of the genetic contributions to cancer and other adult diseases. Our study includes publicly posted genomic summary results from the PLCO Atlas<sup>61</sup>. No institutional review board review is required for PLCO summary data use. The Vanderbilt University Medical Center institutional review board approved the BioVU study. The Malmö Diet and Cancer Study (MDCS) was approved by the local ethics committee.

### Study populations and phenotyping

Genome-wide association analyses of PSA levels were conducted using germline genetic data derived from DNA extracted from non-prostatic tissues (for example, blood and buccal swabs). Analyses were restricted to *cis*-gender men, defined as individuals of biological male sex and self-reported male gender identity who had never been diagnosed with prostate cancer. Men with a history of surgical resections of the prostate were excluded in studies for which this information was available. To reduce potential for reverse causation, analyses were limited to PSA values  $\leq 10$  ng ml<sup>-1</sup>, which corresponds to low-risk prostate cancer based on the D'Amico prostate cancer risk classification system<sup>62</sup>, and PSA  $> 0.01$  ng ml<sup>-1</sup>, to ensure that individuals had a functional prostate not impacted by surgery or radiation.

The UKB is a population-based prospective cohort of over 500,000 individuals aged 40–69 years at enrollment in 2006–2010 with genetic and phenotypic data<sup>63</sup>. Health-related outcomes were ascertained via individual record linkage to national cancer and mortality registries and hospital inpatient encounters. PSA values were abstracted from primary care records for a subset of participants with genetic data. Field code mappings used to identify PSA values included any serum PSA measure except for free PSA or ratio of free to total PSA (Supplementary Table 25).

The Kaiser Permanente GERA cohort used in this analysis was previously described in Hoffmann et al.<sup>17</sup>. In brief, prostate cancer status was ascertained from the Kaiser Permanente Northern California Cancer Registry, the Kaiser Permanente Southern California Cancer Registry or through review of clinical electronic health records. PSA levels were abstracted from Kaiser Permanente electronic health records from 1981 through 2015.

The PLCO Cancer Screening Trial is a completed randomized trial that enrolled approximately 155,000 participants between November 1993 and July 2001. The PLCO Cancer Screening Trial was designed to determine the effects of screening on cancer-related mortality and secondary endpoints in men and women aged 55–74 years<sup>64</sup>. Men randomized to the screening arm of the trial underwent annual screening with PSA for 6 years and digital rectal exam (DRE) for 4 years<sup>64</sup>. These analyses were limited to men with a baseline PSA measurement who were randomized to the screening arm of the trial ( $n = 29,524$ ). Men taking finasteride at the time of PSA measurement were excluded from analysis.

The Vanderbilt University Medical Center BioVU resource is a synthetic derivative biobank linked to de-identified electronic health records<sup>65</sup>. Analyses were based on PSA levels that were measured as part of routine clinical care.

The MDCS is a population-based prospective cohort study that recruited men and women aged between 44 years and 74 years of age who were living in Malmö, Sweden between 1991 and 1996 to investigate the impact of diet on cancer risk and mortality<sup>66</sup>. These analyses included men from the MDCS who were not diagnosed with prostate cancer as of December 2014 and had available genotyping and baseline PSA measurements<sup>66</sup>.

The PCPT is a completed phase 3 randomized, double-blind, placebo-controlled trial of finasteride for prostate cancer prevention that began in 1993 (ref. 8). The PCPT randomly assigned 18,880 men aged 55 years or older who had a normal DRE and PSA level  $\leq 3$  ng ml<sup>-1</sup> to either finasteride or placebo. For men with multiple pre-randomization PSA values, the earliest value was selected. Cases included all histologically confirmed prostate cancers detected during the 7-year treatment period and tumors that were detected by the end-of-study prostate biopsy. Our analyses included the subset of PCPT participants that was genotyped on the Illumina Infinium Global Screening Array 24 v2.0.

The SELECT is a completed phase 3 randomized, placebo-controlled trial of selenium (200 µg per day from L-selenomethionine) and/or vitamin E (400 IU per day of all *rac*- $\alpha$ -tocopheryl acetate) supplementation for prostate cancer prevention<sup>37</sup>. Between 2001 and 2004, 34,888 eligible participants were randomized. The minimum enrollment age was 50 years for African American men and 55 years for all other men<sup>37</sup>. Additional eligibility requirements included no prior prostate cancer diagnosis,  $\leq 4$  ng ml<sup>-1</sup> of PSA in serum and a DRE not suspicious for cancer. For men who had multiple pre-randomization PSA values, the earliest value was selected. Our analyses included a subset of SELECT participants genotyped on the Illumina Infinium Global Screening Array 24 v2.0.

### Quality control and genome-wide association analyses

Standard genotyping and quality control (QC) procedures were implemented in each participating study. Before meta-analysis, we applied variant-level QC filters that included low imputation quality (INFO  $< 0.30$ ), MAF  $< 0.005$  and deviations from Hardy–Weinberg equilibrium ( $P_{\text{HWE}} < 1 \times 10^{-5}$ ). Sample-level filtering was performed to remove samples with discordant genetic sex and self-reported gender and call rate  $< 0.97$ . One sample from each pair of first-degree relatives was also excluded. GWAS phenotypes and adjustment covariates are reported in Supplementary Table 26. Genome-wide association analyses performed linear regression of log(PSA) as the outcome, using age and genetic ancestry principal components (PCs) as the minimum set of covariates.

**UKB.** Genotyping and imputation for the UKB cohort were previously described<sup>63</sup>. In brief, participants were genotyped on the UKB Affymetrix Axiom array (89%) or the UK BiLEVE array (11%) with imputation performed using the Haplotype Reference Consortium (HRC) and the merged UK10K and 1000 Genomes phase 3 reference panels. Genetic ancestry PCs were computed using fastPCA based on a set of 407,219 unrelated samples and 147,604 genetic markers<sup>63</sup>. Association analyses in the UKB were restricted to individuals of European ancestry based on self-report ('White') and after excluding samples with either of the first two genetic ancestry PCs outside of 5 s.d. of the population mean, as previously described<sup>49</sup>. We removed samples with discordant self-reported and genetic sex as well as one sample from each pair of first-degree relatives identified using KING<sup>67</sup>. Using a subset of genotyped autosomal variants with MAF  $\geq 0.01$  and call rate  $\geq 97\%$ , we filtered samples with heterozygosity  $> 5$  s.d. from the mean. For participants with multiple PSA measurements, the median value PSA was used. Sensitivity analyses were conducted comparing this approach to a GWAS of individual-specific random effects derived from fitting a linear mixed model to repeated log(PSA) values.

**GERA.** Genotyping, imputation and QC of the GERA cohort were previously described<sup>17,68,69</sup>. In brief, all men were genotyped for over 650,000 single-nucleotide polymorphisms (SNPs) on four race/ethnicity-specific Affymetrix Axiom arrays that were optimized for individuals who self-identified as non-Hispanic white, Latino, East Asian and African American, respectively<sup>68,69</sup>. Genotype QC procedures and imputation for the original GERA cohort were performed on an array-wise basis, as previously described<sup>17,70</sup>. Pre-phasing was done by SHAPEIT version 2.5 (ref. 71) and imputation with IMPUTE2 version 2.3.1 (ref. 72) using the 1000 Genomes phase 3 release with 2,504 samples. The top ten genetic ancestry PCs from EIGENSTRAT version 4.2 were included in the linear model as ancestry covariates<sup>73</sup>. Analyses were conducted according to self-identified race/ethnicity groups. Residuals were computed from linear mixed models that were fit to repeated log(PSA) measures. This approach was nearly identical to a long-term average, except that it used the median instead of the mean to handle any potential outlier PSA level values.

**PLCO Atlas.** Our study used GWAS summary statistics from the PLCO Atlas Project, a resource for multi-trait GWAS. Genotyping, QC and imputation procedures for this resource are described by Machiela et al.<sup>61</sup>. The Atlas Project combined genotyping data previously generated by high-density arrays for 25,831 participants (OncoArray, Omni2.5M and OmniExpress) with a new round of genotyping using the Illumina Global Screening Array (GSA). For participants genotyped on multiple genotyping arrays ( $n = 1,192$ ), data from only one array were retained, with the following prioritization: GSA > OncoArray > Omni2.5M > OmniExpress. Extensive QC filtering was performed for subsequent imputation and association analyses. Iterative 80% and 95% sample-level and variant-level call rate filters were applied to remove poorly genotyped samples and variants. Samples with > 20% estimated contamination based on VerifyIDintensity<sup>74</sup> were also removed. Samples with discordant self-reported gender and genetically inferred sex were identified based on X-chromosome method-of-moments  $F$  coefficient from PLINK, using 0.5 as the threshold ( $F$  coefficients are close to 0.0 for males and 1.0 for females). Heterozygosity outliers were detected using absolute values from PLINK method-of-moments  $F$  coefficients > 0.2.

Genetic ancestry was determined using GRAF<sup>75</sup> on a set of 10,000 pre-selected fingerprinting variants. Participants were assigned to nine ancestral groups: 'African', 'African American', 'East Asian', 'European', 'Hispanic1', 'Hispanic2', 'Other', 'Other Asian' and 'South Asian'. Hispanic1 included individuals of Dominican or Puerto Rican ancestry, whereas Hispanic2 included individuals of Mexican or Latin American ancestry. For parsimony, we merged 'African' and 'African American' into an 'African American (Combined)' and 'East Asian' and 'Other Asian' into an 'East Asian (Combined)'. Imputation was performed using the TOPMed 5b reference panel, which is accessible via the TOPMed Imputation Server hosted on the Michigan Imputation Server. Before imputation, variants with  $MAF \leq 0.01$ , missingness  $\geq 0.05$  and Hardy–Weinberg deviations ( $P_{HWE} \leq 1 \times 10^{-6}$ ) were removed. Genotyped data were aligned to reference datasets using a community-recommended script (HRC-1000G-check-bim.pl from <https://www.well.ox.ac.uk/~wrayner/tools/>) that was modified to support the TOPMed 5b reference panel using a pre-existing test imputation with 1000 Genomes subjects. Pre-phasing using phased reference data from TOPMed release 5b was conducted using Eagle 2.4 (ref. 76). Imputation was conducted against the same reference panel using minimac4. GWAS was based on the first PSA value for each PLCO participant.

**BioVU.** Participants were identified using Vanderbilt University Medical Center's BioVU resource, a DNA biobank comprising ~270,000 individuals and linked to a de-identified electronic health record<sup>65</sup>. All participants ( $n = 8,074$ ) were genotyped on Illumina's Expanded Multi-Ethnic Genotyping Array (MEGA<sup>EX</sup>) platform. Genetic ancestries

were assigned by running principal component analysis using SNPRelate<sup>77</sup> on a set of pruned SNPs ( $Rsq < 0.5$ ,  $MAF \geq 0.1$ ). Participants were classified as European ancestry if their first two PCs were within 4 s.d. of the median for the participants reporting 'White' as their race. Participants were classified as African ancestry if their first two PCs were within 4 s.d. of the median for participants reporting their race as 'Black'. All QC procedures were performed using PLINK version 1.90. We removed one randomly selected sample out of each pair of related individuals ( $\pi\text{-hat} \geq 0.2$ ) identified using identity-by-descent. We excluded participants with SNP missingness > 3% or heterozygosity > 5 s.d. from the mean. Before imputation, data were pre-processed using the HRC-1000G-check-bim.pl (from <http://www.well.ox.ac.uk/~wrayner/tools/>) and pre-phased using Eagle version 2.4 (ref. 76). Genetic data were imputed on the Michigan Imputation Server using 1000 Genomes phase 3 version 5 as the reference panel. For men with multiple PSA measurements, the median PSA was used.

**MDCS.** Data from multiple batches of genotyping of 4,069 MDCS participants using different Illumina Omni arrays were merged. For variants that appeared more than once under different names on the same Illumina array, those with the higher genotyping rate were retained. Indels, ambiguous palindromic (for example, A/T or C/G alleles) and multi-allelic variants were removed. Only SNPs that we could unambiguously map to the 1000 Genomes phase 1 dataset were kept. Individuals with > 10% missingness were removed. Next, SNPs with a missingness rate > 10% or deviation from Hardy–Weinberg equilibrium ( $P_{HWE} < 0.001$ ) were removed. At this stage, the PCs of ancestry were computed. Individuals for whom the inferred sex based on X-chromosome heterozygosity was not male, or for whom there were more than two genetic mismatches with 40 SNPs that we had previously genotyped in these samples with targeted genotyping<sup>66</sup>, were excluded.

To assess genetic ancestry, MDCS data were combined with data from HapMap phase 3 for variants present in all genotyping batches. These SNPs were further filtered to have < 0.01% missingness and LD pruned ( $-indep\text{-}parwise\ 50\ 50.05$ ). SMARTPCA in EIGENSOFT (<https://github.com/chrchang/eigensoft>) was run on the resulting 18,299 SNPs to generate the top ten genetic ancestry PCs. Analyses were restricted to individuals of European ancestry based on clustering with HapMap reference populations and exclusion of outliers with a  $z$ -score on PC1 and PC2 > 5. Imputation was performed using the TOPMed 5b reference panel, which is accessible via the TOPMed Imputation Server hosted on the Michigan Imputation Server. Before imputation, the input file was aligned to the build37 reference genome on the basis of chromosome, position and alleles. A total of 847,133 SNPs that passed pre-imputation QC were uploaded to the imputation server. From the resulting imputed files, analyses were restricted to individuals without a prostate cancer diagnosis by 31 December 2014, with individual missingness < 3% and a  $z$ -score < 5.0 for heterozygosity. Log(PSA) values were analyzed using robust linear regression with Tukey biweights. GWAS was performed using linear regression on the residuals extracted from the fitted models.

**PCPT and SELECT.** Participants from PCPT and SELECT were genotyped on the Illumina Infinium Global Screening Array 24 v2.0 and underwent the same QC and imputation procedures. Genotyping calling and QC were performed at the Center for Inherited Disease Research at Johns Hopkins. After removal of samples that failed to produce valid output during initial processing and clustering, the completion rate was 0.9951 and 0.9959 in PCPT and SELECT, respectively. A two-stage filter by completion rate threshold of 0.8 for samples and 0.8 for variants, followed by 0.95 for samples and 0.95 for variants, was performed. Samples with discordant self-reported gender and genetically inferred sex were identified based on X-chromosome method-of-moments  $F$  coefficient from PLINK, using 0.5 as the threshold ( $F$  coefficients are



close to 0.0 for males and 1.0 for females). Identity-by-descent for all subject pairs was determined using PLINK, with close (first and second degree) relatives identified based on a threshold of 0.20. One randomly selected sample from each pair of relatives was retained.

Ancestry was estimated using a set of LD-pruned markers and running SNPWEIGHTS<sup>78</sup> with the reference panel provided containing the following populations: European, West African and East Asian, with a threshold of 0.8 used for imputed ancestry designation. Participants were assigned to a single ancestry group if the ancestry score was  $\geq 0.80$  for just one group. Participants were assigned to an admixed cluster if their ancestry score was  $> 0.20$  and  $< 0.80$  for only one group (for example, ADMIXED\_AFR where AFR = 0.75, EUR = 0.17, EAS = 8). Intermediate ancestry clusters included individuals with ancestry scores matching those criteria in multiple groups:  $0.20 < \text{AFR\_EUR} < 0.80$  (for example, AFR = 0.65, EUR = 0.33) and  $0.20 < \text{EAS\_EUR} < 0.80$  (for example, EUR = 0.55, EAS = 0.43). Autosomal heterozygosity was assessed using the method-of-moments  $F$  coefficient calculated within each ancestry cluster. Heterozygosity outliers were identified and excluded using a threshold of 0.10. Principal component analysis was performed with SMARTPCA in EIGENSOFT (<https://github.com/chrchang/eigensoft>) on a set of LD-pruned markers after splitting by ancestry cluster, to resolve more detailed population substructure. Genetic ancestry PCs were not computed for small clusters ( $n < 50$ ) or individuals who failed other QC filters. For validation of  $\text{PGS}_{\text{PSA}}$  in PCPT and SELECT, we combined ADMIXED\_AFR and AFR\_EUR and treated this as a single group with admixed AFR and EUR ancestry proportions (AFR/EUR). ADMIXED\_EAS and EAS\_EUR were also combined into a single cluster with admixed EAS and EUR ancestry (EAS/EUR).

To prepare genotype data for imputation with the TOPMed 5b reference panel, variants with  $\text{MAF} < 0.001$ , call rate  $< 98\%$  or evidence of deviation from Hardy–Weinberg equilibrium ( $P_{\text{HWE}} < 10^{-6}$ ) were removed. After these QC steps, a total of 474,046 variants remained for PCPT, and 491,015 variants were retained for SELECT. Before submitting the data to the TOPMed Imputation Server, files were pre-processed using the check-bim.pl script (<http://www.well.ox.ac.uk/~wrayner/tools/>). Next, chromosomal positions were lifted over from GRCh37/hg19 to GRCh38 and aligned against the TOPMed reference SNP list based on chromosome, position and alleles to ensure that reference and alternate alleles were correct in the resulting VCF files.

### Heritability of PSA levels attributed to common variants

Heritability of PSA levels was estimated using individual-level data and GWAS summary statistics. UKB participants with available PSA and genetic data were analyzed using LDK version 5.1 (ref. 24) and GCTA version 1.93 (ref. 23), following the approach previously implemented in the GERA cohort<sup>17</sup>. Genetic relationship matrices were filtered to ensure that no pairwise relationships with kinship estimates  $> 0.05$  remained. Heritability was estimated using common ( $\text{MAF} \geq 0.01$ ) LD-pruned ( $r^2 < 0.80$ ) variants with imputation INFO  $> 0.80$ . We implemented the LDK-Thin model using the recommended genetic relatedness matrix (GRM) settings (INFO  $> 0.95$ , LD  $r^2 < 0.98$  within 100 kb) and the same parameters as GCTA for comparison (LD  $r^2 < 0.80$ , INFO  $> 0.80$ ). For both methods, sensitivity analyses were conducted using more stringent GRM settings (kinship = 0.025, genotyped variants).

Summary statistics from GWAS results based on the same set of UKB participants ( $n = 26,491$ ) and from a European ancestry GWAS meta-analysis ( $n = 85,824$ ) were analyzed using LDK, LD score regression (LDSR)<sup>25</sup> and an extension of LDSR using a high-definition likelihood (HDL) approach<sup>26</sup>. For LDSR, we used the default panel comprising variants available in HapMap3 with weights computed in 1000 Genomes version 3 EUR individuals and in-house LD scores computed in UKB European ancestry participants<sup>49</sup>. The baseline linkage disequilibrium (BLD)-LDK model was fit using pre-computed tagging files calculated in UKB GBR (white British) individuals for HapMap3 variants from the LDSR default panel. HDL analyses were conducted using

the UKB-derived panel restricted to high-quality imputed HapMap3 variants<sup>26</sup>. All GWAS summary statistics had sufficient overlap with the reference panels, not exceeding the 1% missingness threshold for HDL and the 5% missingness threshold for LDK and LDSR.

### Genome-wide meta-analysis

Each ancestral population was analyzed separately, and GWAS summary statistics were combined via meta-analysis (Fig. 1). We first used METAL<sup>79</sup> to conduct an inverse-variance-weighted fixed-effects meta-analysis in each ancestry group and then meta-analyzed the ancestry-stratified results. Multi-ancestry meta-analysis results were processed using clumping to identify independent association signals by grouping variants based on LD within specific windows. Clumps were formed around index variants with the lowest genome-wide significant ( $P < 5 \times 10^{-8}$ ) meta-analysis  $P$  value. All other variants with LD  $r^2 > 0.01$  within a  $\pm 10$ -Mb window were considered non-independent and assigned to that lead variant. Since over 90% of the meta-analysis consisted of individuals of European ancestry, clumping was performed using 1000 Genomes phase 3 EUR and UKB reference panels, which yielded concordant results. We confirmed that LD among the resulting lead variants did not exceed  $r^2 = 0.05$  using a merged 1000 Genomes ALL reference panel.

We first examined heterogeneity in the multi-ancestry fixed-effects meta-analysis results using Cochran's  $Q$  statistic. To assess heterogeneity specifically due to ancestry, we applied MR-MEGA<sup>27</sup>, a meta-regression approach for aggregating GWAS results across diverse populations. Summary statistics from each GWAS were meta-analyzed using MR-MEGA without combining by ancestry first. The MR-MEGA analysis was performed across four axes of genetic variation derived from pairwise allele frequency differences, based on the recommendation for separating major global ancestry groups. Index variants from the MR-MEGA analysis were selected using the same clumping parameters as described above (LD  $r^2 < 0.01$  within a  $\pm 10$ -Mb window), based on the merged 1000 Genomes ALL reference panel. For each variant, we report two heterogeneity  $P$  values: one that is correlated with ancestry and accounted for in the meta-regression ( $P_{\text{Het-Anc}}$ ) and the residual heterogeneity that is not due to population genetic differences ( $P_{\text{Het-Res}}$ ).

### $\text{PGS}_{\text{PSA}}$ development and validation

We implemented two strategies for generating a genetic score for PSA levels. In the first approach, we selected 128 variants that were genome-wide significant ( $P < 5 \times 10^{-8}$ ) in the multi-ancestry meta-analysis and were independent (LD  $r^2 < 0.01$  within a  $\pm 10$ -Mb window) in 1000 Genomes EUR and (LD  $r^2 < 0.05$ ) 1000 Genomes ALL populations ( $\text{PGS}_{128}$ ). Each variant in  $\text{PGS}_{128}$  was weighted by the meta-analysis effect size estimated using METAL. As an alternative strategy to clumping and thresholding, we fit a genome-wide score using the PRS-CSx algorithm<sup>36</sup>, which takes GWAS summary statistics from each ancestry group as inputs and estimates posterior SNP effect sizes under coupled continuous shrinkage priors across populations ( $\text{PGS}_{\text{CSx}}$ ). Analyses were conducted using pre-computed population-specific LD reference panels from the UKB, which included 1,287,078 HapMap3 variants that are available in both the UKB and 1000 Genomes phase 3.

We calculated a single trans-ancestry PGS that can be applied to all participants in the target cohort, rather than optimizing a PGS within each ancestry group. This approach is more robust to differences in genetic ancestry assignments across studies and does not require separate testing and validation datasets for parameter tuning each ancestry group<sup>36</sup>. To facilitate this type of analysis, PRS-CSx provides a –meta option that integrates population-specific posterior SNP effects using an inverse-variance-weighted meta-analysis in the Gibbs sampler<sup>36</sup>. The global shrinkage parameter was set to  $\phi = 0.0001$ . PRS-CSx was run on the intersection of variants that were in the LD reference panel and had imputation quality (INFO  $> 0.90$ ), resulting in 1,058,163 variants in PCPT and 1,071,268 variants in SELECT. Because PRS-CSx considers only

autosomes, chrX variants that were included in PGS<sub>128</sub> were added to PGS<sub>CSX</sub> separately, when output files from each chromosome produced by the PLINK–score command were concatenated.

The predictive performance of PGS<sub>CSX</sub> and PGS<sub>128</sub> was evaluated in two independent cancer prevention trials that were not included in the meta-analysis: PCPT and SELECT. Analyses were conducted in the pooled sample for each cohort, which included individuals of all ancestries who passed QC filters (Supplementary Note). Ancestry-stratified analyses were conducted for clusters with  $n > 50$  with available genetic ancestry PCs. Ancestry scores were computed with SNPWEIGHTS<sup>78</sup>. Individuals with ancestry scores  $\geq 0.80$  for a single group were assigned to clusters for predominantly European (EUR), West African (AFR) and East Asian (EAS) ancestry. Admixed individuals with intermediate ancestry scores for at least one group were assigned to separate clusters:  $0.20 < \text{EUR/AFR} < 0.80$  or  $0.20 < \text{EUR/EAS} < 0.80$ . Pooled analyses were adjusted for ten within-cluster PCs and global ancestry proportions (AFR and EAS).

### Index event bias analysis

Index event bias occurs when individuals are selected based on the occurrence of an event or specific criterion. This is analogous to the direct dependence of one phenotype on another, as in the commonly used example of cancer survival<sup>34</sup>. Due to unmeasured confounding, this dependence can induce correlations between previously independent risk factors among those selected<sup>33,34</sup>. Genetic effects on prostate cancer can be viewed as conditional on PSA levels, because elevated PSA typically triggers diagnostic investigation. Genetic factors resulting in higher constitutive PSA levels may also increase the likelihood of prostate cancer detection due to more frequent testing (Fig. 4). This selection mechanism could bias prostate cancer GWAS associations by capturing both direct genetic effects on disease risk and selection-induced PSA signals. In the GWAS setting, methods using summary statistics have been developed to estimate and correct for this bias<sup>33,35</sup>. Although typically derived assuming a binary selection trait, these methods are still applicable to selection or adjustment based on quantitative phenotypes<sup>33</sup>. In this study, we conceptualized PSA variation as the selection trait and prostate cancer incidence as the outcome trait (Fig. 4).

We applied the method described in Dudbridge et al.<sup>33</sup>, which tests for index event bias and estimates the corresponding correction factor ( $b$ ) by regressing genetic effects on the selection trait (PSA) against their effects on the subsequent trait (prostate cancer), with inverse variance weights:  $w = 1/(\text{SE}_{\text{PrCa}})^2$ . Summary statistics for prostate cancer were obtained from the most recent prostate cancer GWAS from the PRACTICAL consortium<sup>32</sup>. Sensitivity analyses were performed using Slope-Hunter<sup>35</sup>, an extension of the Dudbridge approach that allows for direct genetic effects on the index trait and subsequent trait to be correlated. For both methods, analyses were conducted using relevant summary statistics and 127,906 variants pruned at the recommended threshold<sup>33</sup> ( $\text{LD } r^2 < 0.10$  in 250-kb windows) with  $\text{MAF} \geq 0.05$  in the 1000 Genomes EUR reference panel. After merging the pruned 1000 Genomes variants with each set of summary statistics, variants with large effects, ( $|\beta| > 0.20$ ) on either  $\log(\text{PSA})$  or prostate cancer, were excluded. The resulting estimate ( $b$ ), adjusted regression dilution using the SIMEX algorithm, was used as a correction factor to recover unbiased genetic effects for each variant:  $\beta'_{\text{PrCa}} = \beta_{\text{PrCa}} - b \times \beta_{\text{PSA}}$ , where  $\beta_{\text{PSA}}$  is the per-allele effect on  $\log(\text{PSA})$ , and  $\beta_{\text{PrCa}}$  is the  $\log(\text{OR})$  for prostate cancer.

The impact of the bias correction was assessed in three ways. First, genome-wide significant prostate cancer index variants were selected from the European ancestry PRACTICAL GWAS meta-analysis (85,554 cases and 91,972 controls) using clumping ( $\text{LD } r^2 < 0.01$  within 10 Mb) (ref. 32). We tabulated the number of variants that remained associated at  $P < 5 \times 10^{-8}$  after bias correction. Next, we fit genetic scores for PSA and prostate cancer in men of European ancestry in the UKB who were not included in the PSA or prostate cancer GWAS (11,568 prostate

cancer cases and 152,884 controls). We compared the correlation between the PGS for PSA (PGS<sub>PSA</sub>), comprising 128 lead variants, and the 269-variant prostate cancer risk score fit with original risk allele weights (PGS<sub>269</sub>) and with weights corrected for index event bias (PGS<sub>269</sub><sup>adj</sup>). To allow adjustment for genetic ancestry PCs and genotyping array, associations between the two scores were estimated using linear regression models. Next, we examined associations for each genetic score (PGS<sub>269</sub>, PGS<sub>269</sub><sup>adj</sup>, PGS<sub>269</sub><sup>adj-S</sup>) with prostate cancer in a subset of GERA participants who underwent a biopsy. Because GERA controls were included in the PSA GWAS meta-analysis, AUC estimates and corresponding bootstrapped 95% CIs were obtained using tenfold cross-validation. We also examined PGS associations with Gleason score, a marker of disease aggressiveness, which was not available in the UKB. Multinomial logistic regression models with Gleason score  $\leq 6$  (reference), 7 and  $\geq 8$  as the outcome were fit for each score in 4,584 cases from the GERA cohort.

### Application of genetically adjusted PSA for biopsy referral and prostate cancer detection

Genetically corrected PSA values were calculated for individual  $i$  as follows<sup>17,19</sup>:

$$\text{PSA}_i^G = \frac{\text{PSA}_i}{a_i} \quad (1)$$

where  $a_i$  is a personalized adjustment factor derived from PGS<sub>PSA</sub>. Because genetic effects were estimated for  $\log(\text{PSA})$ ,  $a_i$  for correcting PSA in  $\text{ng ml}^{-1}$  was derived as:

$$a_i = \frac{\exp(\text{PGS}_i)}{\overline{\exp(\text{PGS})}} \quad (2)$$

$\overline{\text{PGS}}$  can be estimated in controls without prostate cancer or obtained from an external control population<sup>17,19</sup>. We see that  $a_i > 1$  when an individual has a higher multiplicative increase in PSA than the sample average due to their genetic profile, resulting in a lower genetically adjusted PSA compared to the observed value ( $\text{PSA}_i^G < \text{PSA}_i$ ).

We evaluated the potential utility of PGS<sub>PSA</sub> in two clinical contexts. First, we quantified the impact of using  $\text{PSA}_i^G$  on biopsy referrals by examining reclassification at age-specific PSA thresholds used in the Kaiser Permanente health system. Analyses were conducted in GERA participants with information on biopsy date and outcome, comprising prostate cancer cases not included in the PSA GWAS and controls that were part of the PSA GWAS. To use the same normalization factor for both cases and controls while mitigating bias due to control overlap with the PSA discovery GWAS,  $a_i$  for GERA participants was calculated by substituting  $\overline{\text{PGS}}$  from out-of-sample UKB controls ( $n = 152,884$ ). Upward classification resulting in biopsy eligibility occurred when  $\text{PSA}_i^G > \text{PSA}_i \cap \text{PSA}_i^G > \text{ref}$ , where  $\text{ref}$  is the biopsy referral threshold. Downward classification resulting in biopsy ineligibility was defined as:  $\text{PSA}_i^G < \text{PSA}_i \cap \text{PSA}_i^G < \text{ref}$ . Net reclassification (NR) was summarized separately for cases and controls:

$$\text{NR}_{\text{case}} = P(\text{up}|\text{case}) - P(\text{down}|\text{case})$$

$$\text{NR}_{\text{control}} = P(\text{down}|\text{control}) - P(\text{up}|\text{control})$$

This is equivalent to tabulating the proportion of individuals in each biopsy eligibility category:

$$\text{NR}_{\text{case}} = \left( \frac{n_{\text{eligible}}}{n_{\text{case}}} \right) - \left( \frac{n_{\text{ineligible}}}{n_{\text{case}}} \right)$$

$$\text{NR}_{\text{control}} = \left( \frac{n_{\text{ineligible}}}{n_{\text{control}}} \right) - \left( \frac{n_{\text{eligible}}}{n_{\text{control}}} \right)$$



For each NR proportion, 95% CIs were obtained using the normal approximation:

$$NR \pm 1.96 \times \sqrt{\frac{|NR| \times (1 - |NR|)}{n}}$$

Next, we assessed the performance of risk prediction models for prostate cancer overall, aggressive prostate cancer and non-aggressive prostate cancer in the PCPT and the SELECT.

Because both studies were excluded from the PSA GWAS meta-analysis,  $a_i$  and  $PSA_i^G$  for then PCPT and the SELECT were calculated using  $PGS$  observed in each respective study. Consistent with the  $PGS_{PSA}$  validation analysis, pooled analyses included individuals of all ancestries who passed QC filters. To facilitate ancestry-stratified analyses in SELECT, especially for aggressive disease, we combined AFR and AFR/EUR clusters into a single group (AFR pooled) and similarly pooled EAS and EAS/EUR (EAS pooled). Aggressive prostate cancer was defined as Gleason score  $\geq 7$ ,  $PSA \geq 10 \text{ ng ml}^{-1}$ , T3–T4 stage and/or distant or nodal metastases. We compared AUC estimates for logistic regression models using the following predictors, alone and in combination: baseline PSA, genetically adjusted baseline PSA ( $PSA^G$ )  $PGS_{PSA}$ , prostate cancer risk score with original weights ( $PGS_{269}$ ) (ref. 32) and weights corrected for index event bias ( $PGS_{269}^{adj}$ ).

### Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

### Data availability

UK Biobank data are publicly available by request from <https://www.ukbiobank.ac.uk>. To maintain individuals' privacy, data on the GERA cohort are available by application to the Kaiser Permanente Research Bank (<https://researchbank.kaiserpermanente.org/>). All PLCO genotype data are available in the database of Genotypes and Phenotypes (dbGAP) under accession number [phs001286.v2.p2](https://dbgap.ncbi.nlm.nih.gov/ah/study.cgi?study_id=phs001286.v2.p2) ([https://dbgap.ncbi.nlm.nih.gov/ah/study.cgi?study\\_id=phs001286.v2.p2](https://dbgap.ncbi.nlm.nih.gov/ah/study.cgi?study_id=phs001286.v2.p2)). Companion phenotype data can be requested through the NCI Cancer Data Access System (<https://cdas.cancer.gov/plco/>). GWAS summary statistics are available directly from the PLCO Atlas GWAS Explorer website (<https://exploregwas.cancer.gov/plco-atlas/>) as well as accessed directly through API access (<https://exploregwas.cancer.gov/plco-atlas/#/api-access>). Genome-wide summary statistics for the PSA multi-ancestry meta-analysis and ancestry-stratified summary statistics for the development of the genome-wide PSA polygenic score are available from <https://doi.org/10.5281/zenodo.7460134>. Scoring files for fitting PSA polygenic scores are available from the PGS Catalog: <http://www.pgscatalog.org/score/PGS003378/> and <http://www.pgscatalog.org/score/PGS003379/>.

### Code availability

Genome-wide association analyses were conducted using PLINK version 2.0a3LM (<https://www.cog-genomics.org/plink/2.0/>). Fixed-effects inverse-variance-weighted meta-analysis was performed with METAL using SCHEME STDERR ([https://genome.sph.umich.edu/wiki/METAL\\_Documentation](https://genome.sph.umich.edu/wiki/METAL_Documentation)). Weights for the genome-wide polygenic score for PSA were estimated using PRS-CSx (<https://github.com/getian107/PRS-CSx>). Scripts for fitting polygenic scores, performing the index event bias analysis and calculating genetically adjusted PSA values are available at [https://github.com/lkachuri/precision\\_PSA](https://github.com/lkachuri/precision_PSA).

### References

61. Machiela, M. J. et al. GWAS Explorer: an open-source tool to explore, visualize, and access GWAS summary statistics in the PLCO Atlas. *Sci. Data* **10**, 25 (2023).
62. D'Amico, A. V. Risk-based management of prostate cancer. *N. Engl. J. Med.* **365**, 169–171 (2011).

63. Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
64. Andriole, G. L. et al. Mortality results from a randomized prostate-cancer screening trial. *N. Engl. J. Med.* **360**, 1310–1319 (2009).
65. Roden, D. M. et al. Development of a large-scale de-identified DNA biobank to enable personalized medicine. *Clin. Pharmacol. Ther.* **84**, 362–369 (2008).
66. Klein, R. J. et al. Evaluation of multiple risk-associated single nucleotide polymorphisms versus prostate-specific antigen at baseline to predict prostate cancer in unscreened men. *Eur. Urol.* **61**, 471–477 (2012).
67. Manichaikul, A. et al. Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873 (2010).
68. Hoffmann, T. J. et al. Design and coverage of high throughput genotyping arrays optimized for individuals of East Asian, African American, and Latino race/ethnicity using imputation and a novel hybrid SNP selection algorithm. *Genomics* **98**, 422–430 (2011).
69. Hoffmann, T. J. et al. Next generation genome-wide association tool: design and coverage of a high-throughput European-optimized SNP array. *Genomics* **98**, 79–89 (2011).
70. Kvale, M. N. et al. Genotyping informatics and quality control for 100,000 subjects in the Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort. *Genetics* **200**, 1051–1060 (2015).
71. Delaneau, O., Marchini, J. & Zagury, J. F. A linear complexity phasing method for thousands of genomes. *Nat. Methods* **9**, 179–181 (2011).
72. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G. R. Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat. Genet.* **44**, 955–959 (2012).
73. Banda, Y. et al. Characterizing race/ethnicity and genetic ancestry for 100,000 subjects in the Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort. *Genetics* **200**, 1285–1295 (2015).
74. Jun, G. et al. Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am. J. Hum. Genet.* **91**, 839–848 (2012).
75. Jin, Y., Schaffer, A. A., Sherry, S. T. & Feolo, M. Quickly identifying identical and closely related subjects in large databases using genotype data. *PLoS ONE* **12**, e0179106 (2017).
76. Loh, P. R. et al. Reference-based phasing using the Haplotype Reference Consortium panel. *Nat. Genet.* **48**, 1443–1448 (2016).
77. Zheng, X. et al. A high-performance computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* **28**, 3326–3328 (2012).
78. Chen, C. Y. et al. Improved ancestry inference using weights from external reference panels. *Bioinformatics* **29**, 1399–1406 (2013).
79. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).

### Acknowledgements

The Precision PSA study is supported by funding from the National Institutes of Health (NIH) National Cancer Institute (NCI) under awards R01CA241410 (J.S.W.), U01CA261339 (J.S.W. and D.V.C.) and R00CA246076 (L.K.), and the Young Investigator Award from the Prostate Cancer Foundation (R.E.G.). Contributing studies were supported by research grants from the NIH National Institute of General Medical Sciences (NIGMS) under award R01GM130791 (J.D.M.); NIH/NCI Cancer Center Support Grant to Memorial Sloan Kettering Cancer Center (P30CA008748); MSKCC Specialized Programs of Research Excellence in Prostate Cancer (P50CA92629, H.L.); the Swedish Cancer Society (Cancerfonden 20 1354 PjF, H.L.); and the General Hospital in Malmö Foundation for Combating

Cancer. This work was supported, in part, through the computational resources and staff expertise provided by Scientific Computing at the Icahn School of Medicine at Mount Sinai. Research reported in this paper was supported by the Office of Research Infrastructure of the NIH under award S10OD026880 and NIH/NCI funding (R01CA175491 and R01CA244948, R.J.K.). The content is solely the responsibility of the authors and does not necessarily represent the official views of the NIH.

## Author contributions

Concept and design: L.K., T.J.H., R.J.K., J.D.M., R.E.G. and J.S.W. Acquisition, analysis or interpretation of data: L.K., T.J.H., Y.J., S.I.B., J.P.S., K.S., M.J.M., N.D.F., W.-Y.H., S.A.L., R.E., P.J.G., C.T., I.T., H.L., S.K.V.D.E., S.J.C., C.A.H., D.V.C., R.J.K., J.D.M., R.E.G. and J.S.W. Drafting of the manuscript: L.K., T.J.H., R.E.G. and J.S.W. Critical revision of the manuscript for important intellectual content: L.K., T.J.H., Y.J., S.I.B., J.P.S., K.S., M.J.M., N.D.F., W.-Y.H., S.A.L., R.E., P.J.G., C.T., I.T., H.L., S.K.V.D.E., S.J.C., C.A.H., D.V.C., R.J.K., J.D.M., R.E.G. and J.S.W.

## Competing interests

J.S.W. is a non-employee cofounder of Avail Bio. H.L. is named on a patent for intact PSA assays and a patent for a statistical method to

detect prostate cancer that is licensed to and commercialized by OPKO Health. H.L. receives royalties from sales of the test and has stock in OPKO Health. All other authors have no competing interests.

## Additional information

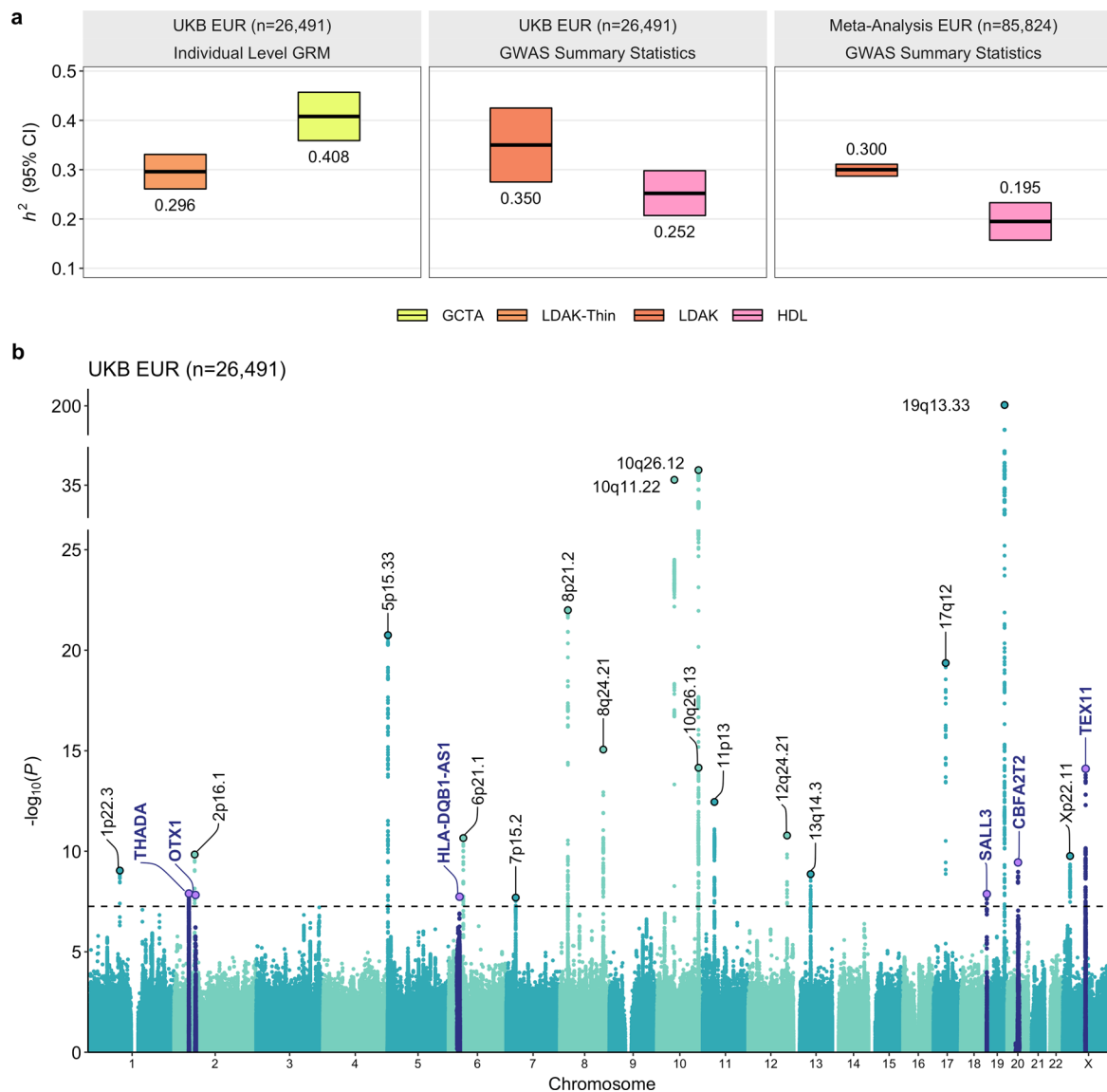
**Extended data** is available for this paper at <https://doi.org/10.1038/s41591-023-02277-9>.

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41591-023-02277-9>.

**Correspondence and requests for materials** should be addressed to Rebecca E. Graff or John S. Witte.

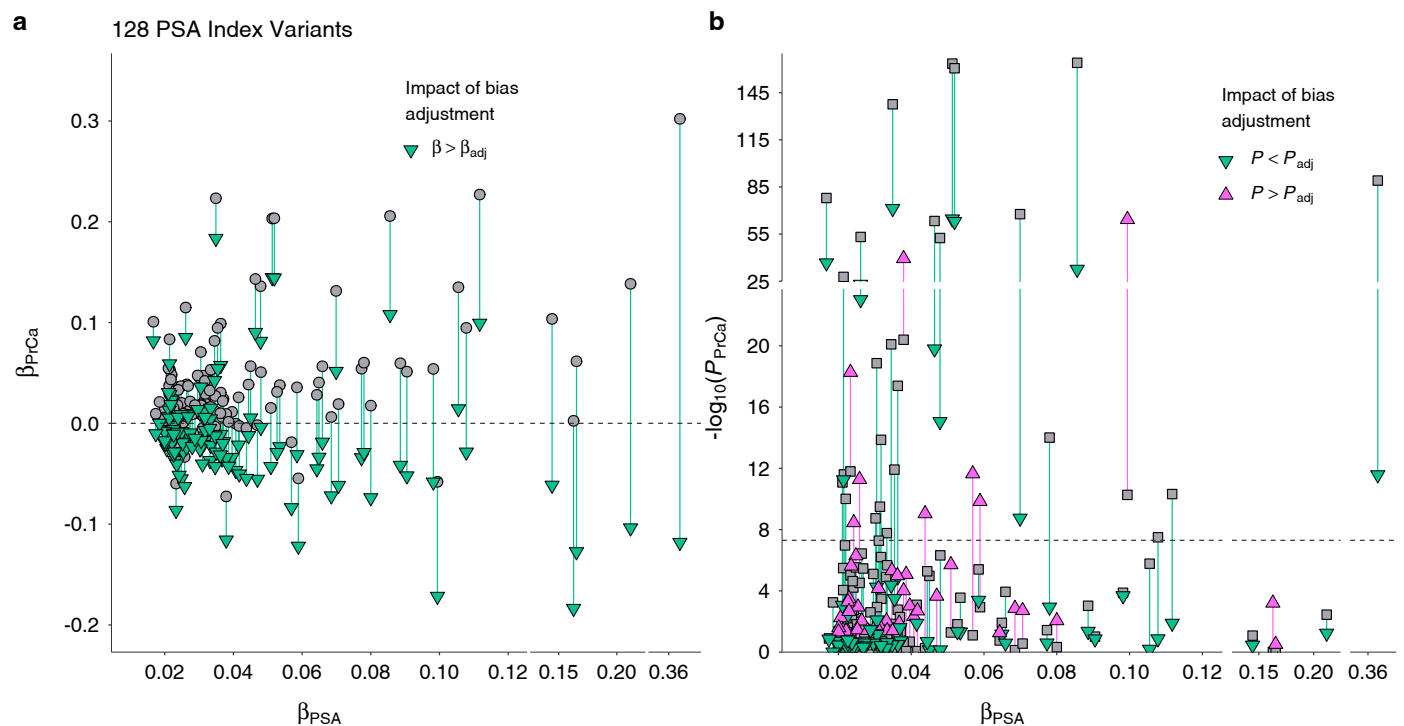
**Peer review information** *Nature Medicine* thanks Jason Vassy, Jian Yang and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Primary handling editor: Anna Maria Ranzoni, in collaboration with the *Nature Medicine* team.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).



**Extended Data Fig. 1 | Heritability ( $h^2$ ) of PSA levels and GWAS results in men of European ancestry without prostate cancer. **a**, Crossbars show  $h^2$  estimates, annotated below, and corresponding 95% confidence intervals across statistical methods. In the UK Biobank (UKB), heritability was estimated using GCTA and Linkage Disequilibrium Adjusted Kinships (LDKA)-Thin models from a genetic relatedness matrix (GRM) of common ( $MAF \geq 0.01$ ) LD-pruned ( $r^2 < 0.80$ ) variants with imputation quality INFO > 0.80. These estimates were compared to analyses**

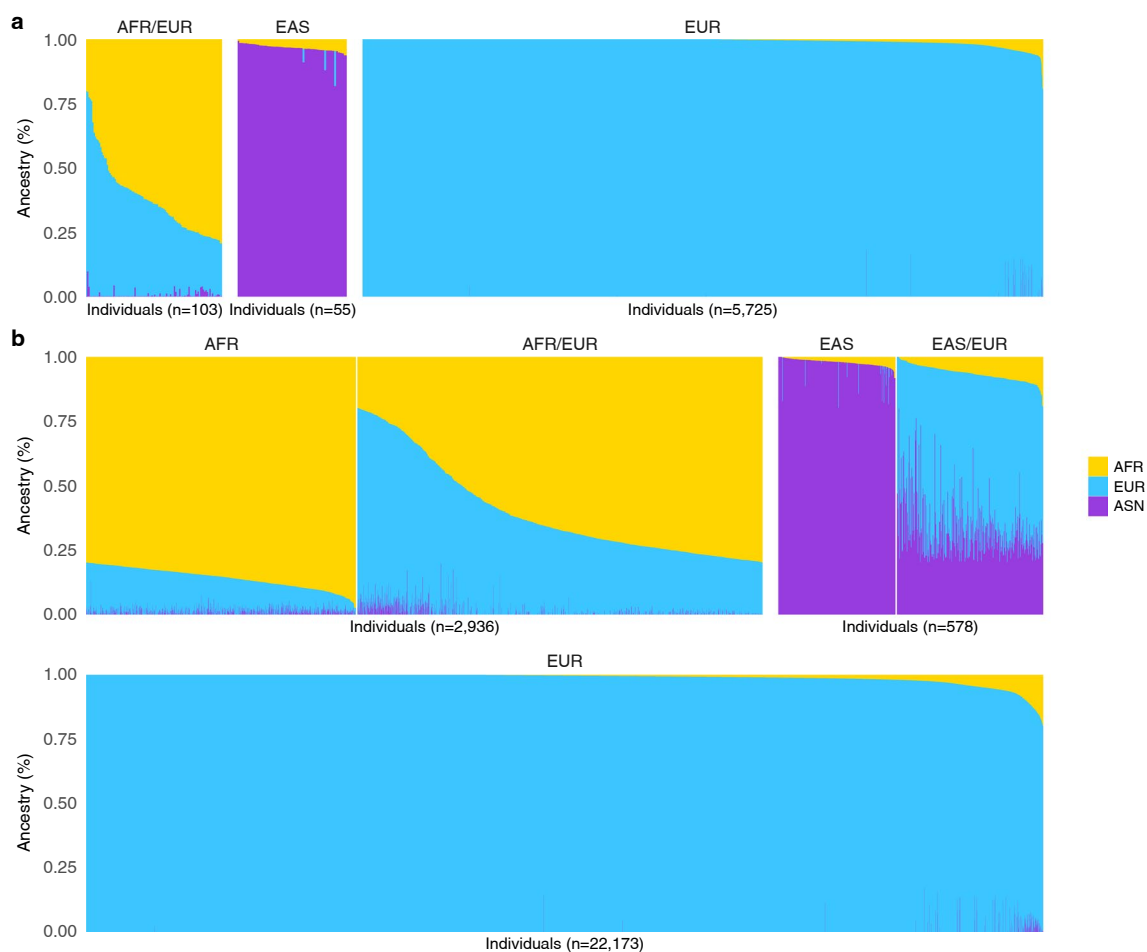
of GWAS summary statistics from the UK Biobank and the EUR meta-analysis using the baseline linkage disequilibrium LDKA model and a high-definition likelihood (HDL) method by Ning et al.<sup>26</sup> **b**, UKB GWAS results where known PSA loci are labeled with the corresponding cytoband region and new regions are labeled with the nearest gene. Highlighted peaks include variants in LD ( $r^2 \geq 0.01$ ) with the lead novel variant. Two-sided p-values are derived from linear regression models.



**Extended Data Fig. 2 | Impact of correction for PSA-related selection bias on genetic associations with prostate cancer.** Associations with prostate cancer for 128 PSA-associated index variants were obtained from the PRACTICAL GWAS by Conti et al.<sup>32</sup> PSA index variants were selected from the multi-ancestry GWAS meta-analysis using clumping and thresholding ( $P < 5 \times 10^{-8}$ , linkage disequilibrium  $r^2 < 0.01$ ). **a**, GWAS effect sizes for prostate cancer ( $\beta_{\text{PrCa}}$ ) are

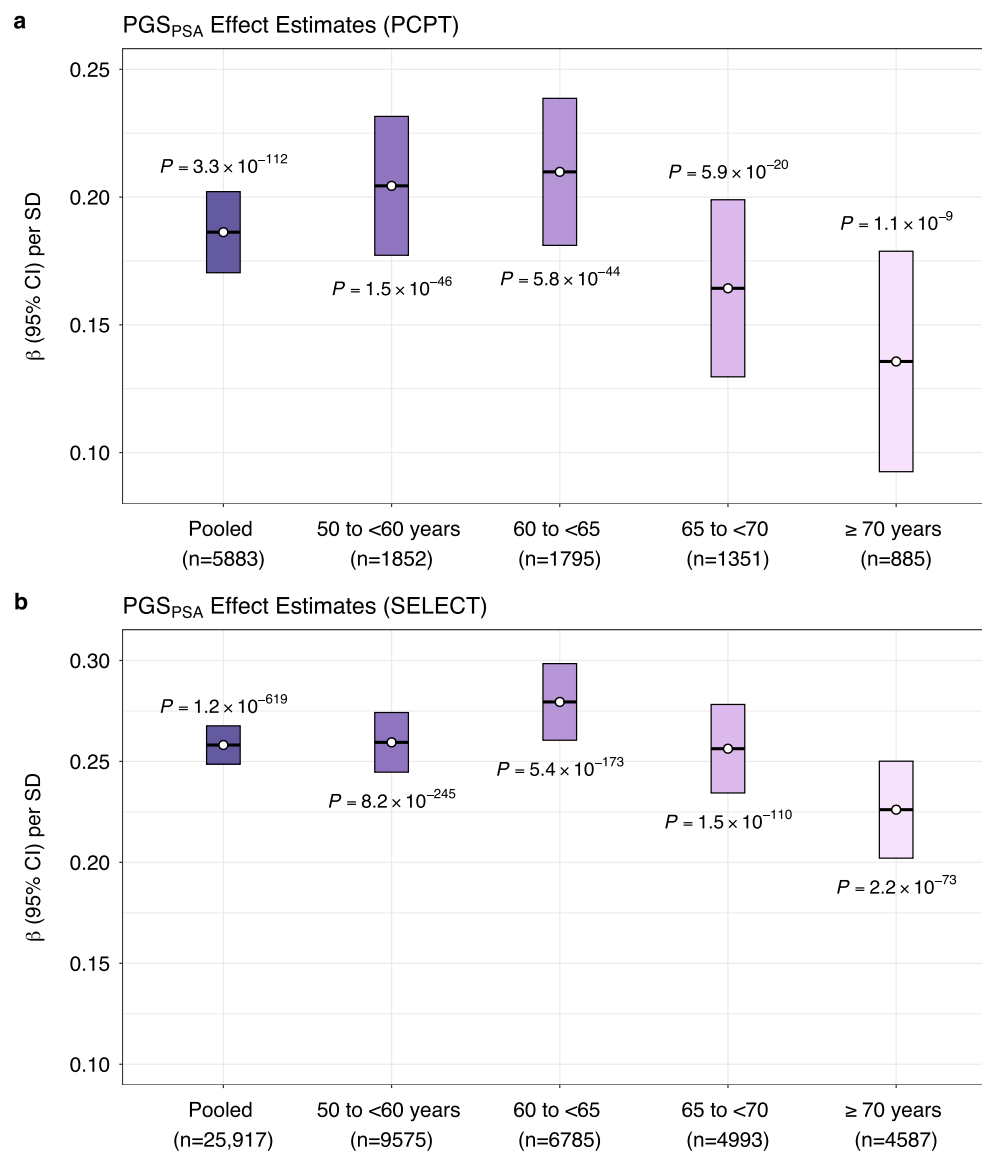
aligned to the PSA-increasing allele. Bias-adjusted effect sizes ( $\beta_{\text{adj}}$ ) are denoted by triangles. **b**, Two-sided GWAS p-values for prostate cancer ( $P_{\text{PrCa}}$ ) were derived from an inverse-variance-weighted fixed-effects meta-analysis. Two-sided bias-adjusted p-values ( $P_{\text{adj}}$ ), denoted by triangles, were calculated from a chi-squared test statistic based on  $\beta_{\text{adj}}$  and corresponding standard errors. Genome-wide significance threshold ( $P < 5 \times 10^{-8}$ ) is indicated by the dotted line.





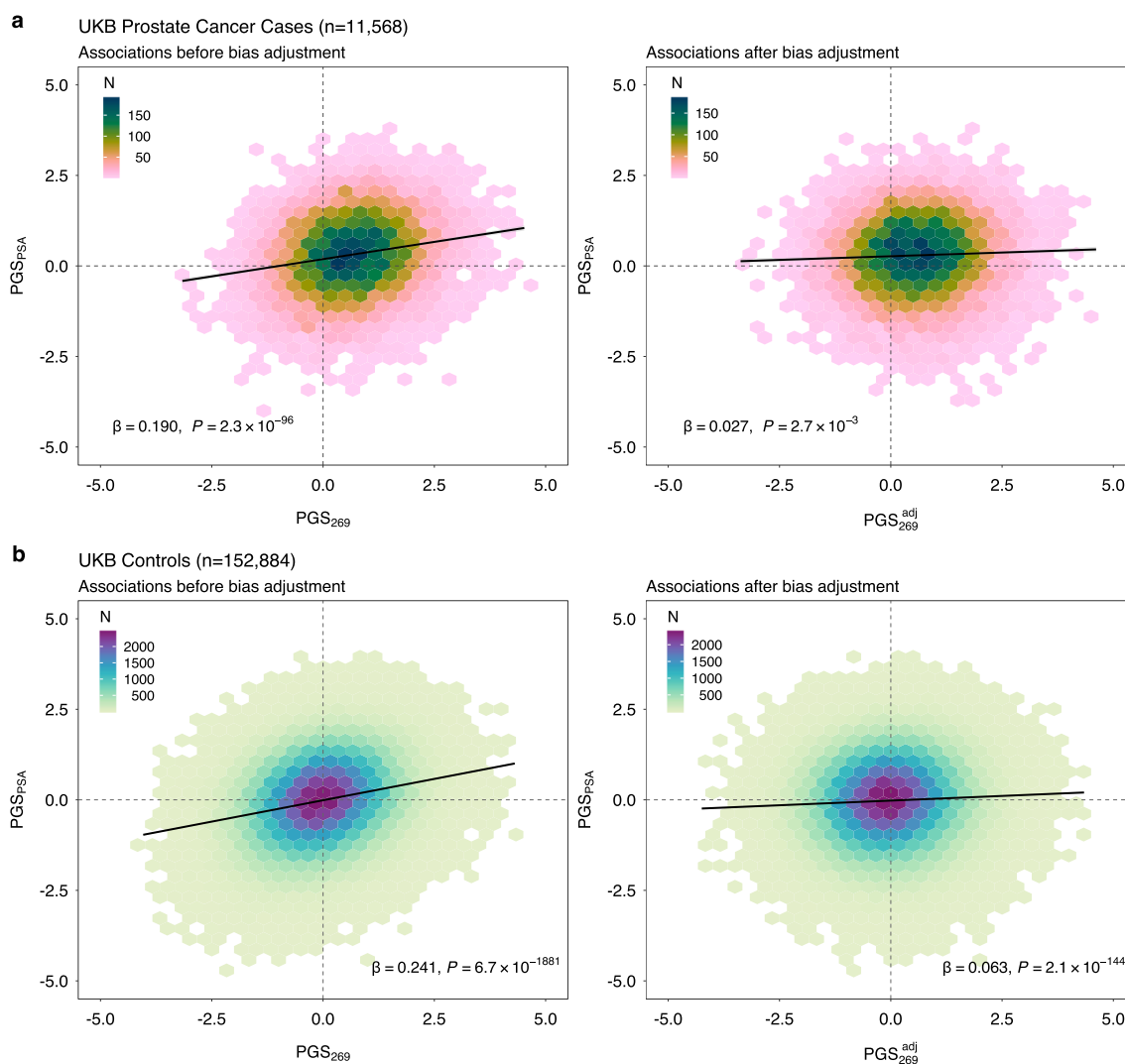
**Extended Data Fig. 3 | Ancestry composition of validation cohorts.** Admixture plots visualizing genetic ancestry proportions for participants within population clusters in **a**, Prostate Cancer Prevention Trial (PCPT) and **b**, Selenium and Vitamin E Cancer Prevention Trial (SELECT). Both cohorts were excluded from the PSA GWAS used for polygenic score development. For each individual, the proportion of African (AFR), European (EUR), and East Asian (EAS) genetic

ancestry is shown. Single-ancestry clusters include individuals with ancestry scores  $\geq 0.80$  in one ancestry group. Admixed ancestry clusters AFR/EUR and EAS/EUR include individuals with ancestry proportions  $>0.20$  and  $<0.80$ . For analyses of prostate cancer risk in SELECT, AFR and AFR/EUR and EAS and EAS/EUR were combined into pooled African ancestry ( $n = 2,936$ ) and pooled East Asian ancestry ( $n = 578$ ), respectively.



**Extended Data Fig. 4 | Age-stratified PGS<sub>PSA</sub> associations.** Performance of the genome-wide PGS<sub>PSA</sub> developed using the PRS-CSx algorithm was evaluated in the two cancer prevention trials: **a**, Prostate Cancer Prevention Trial (PCPT) and **b**, Selenium and Vitamin E Cancer Prevention Trial (SELECT). Crossbars visualize the effect estimates ( $\beta$ ) and corresponding 95% confidence intervals per standard deviation (SD) increase in the standardized PGS<sub>PSA</sub>. Associations between PGS<sub>PSA</sub>

and baseline log(PSA) were estimated in the pooled sample and stratified by age group. All p-values are two-sided and derived from linear regression models adjusted for age at baseline, top 10 population-specific genetic ancestry principal components, and proportions of African and East Asian genetic ancestry.



**Extended Data Fig. 5 | Impact of index event bias on polygenic score (PGS) associations.** Association between PGS for PSA ( $\text{PGS}_{\text{PSA}}$ ) and PGS for prostate cancer ( $\text{PGS}_{269}$ ) fit using original weights, as reported in Conti et al.<sup>32</sup>, is compared to  $\text{PGS}_{269}$  fit using weights that have been adjusted for index event bias ( $\text{PGS}_{269}^{\text{adj}}$ ) using the Dudbridge et al.<sup>33</sup> method. Linear regression lines with shaded 95%

confidence intervals visualizing the PGS associations in **a**, prostate cancer cases and **b**, men not diagnosed with prostate cancer (controls) are overlaid on individual data points summarized as hexbins. Analyses were restricted to male UK Biobank participants of European ancestry who were excluded from the GWAS of PSA levels.

## Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- ☐ ☒ The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- ☐ ☒ A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- ☐ ☒ The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- ☐ ☒ A description of all covariates tested
- ☐ ☒ A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- ☐ ☒ A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- ☐ ☒ For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- ☒ ☐ For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- ☒ ☐ For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- ☐ ☒ Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

### Software and code

Policy information about [availability of computer code](#)

Data collection No software was used for data collection.

Data analysis

- GWAS was performed using PLINK 2.0 (version 2.00a3LM)
- Fixed-effects inverse-variance-weighted meta-analysis was performed with METAL (version 2011), available from: <http://csg.sph.umich.edu/abecasis/Metal/download/>
- Heritability analyses were performed using LDAK (version 5.1), GCTA (version 1.93.2beta), and HDL (version 1.4.0), available from: <https://github.com/zhenin/HDL>
- Other statistical analyses and data visualizations were performed in R (version 4.1.2), including the use of the following R packages:
  - \* SlopeHunter (version 0.0.2), available from: <https://github.com/Osmahmoud/SlopeHunter>
  - \* Polygenic risk score modeling was performed using the PRS-CSx algorithm and reference panels, available from: <https://github.com/getian107/PRScsx> (version 1.0.0. July 29, 2021)
- Scripts for fitting polygenic scores, performing the index event bias analysis, and calculating genetically adjusted PSA values are available from: [https://github.com/lkachuri/precision\\_PSA](https://github.com/lkachuri/precision_PSA)

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.



## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

- The research was conducted with approved access to UK Biobank data under application number 14105 (PI: Witte). UK Biobank data are publicly available by request from <https://www.ukbiobank.ac.uk>.

- To maintain individuals' privacy, data on the GERA cohort are available by application to the Kaiser Permanente Research Bank ([researchbank.kaiserpermanente.org](https://researchbank.kaiserpermanente.org)).

- All PLCO genotype data is available in dbGaP18 under accession number phs001286.v2.p2 (<https://identifiers.org/dbgap:phs001286.v2.p2>). Companion phenotype data can be requested through the NCI Cancer Data Access System (CDAS) (<https://cdas.cancer.gov/plco/>). GWAS summary statistics are available directly from the PLCO Atlas GWAS Explorer website (<https://exploregwas.cancer.gov/plco-atlas/>) as well as accessed directly through API access (<https://exploregwas.cancer.gov/plco-atlas/#/api-access>).

- Scoring files for fitting PSA polygenic scores are available from the PGS Catalog: [www.pgscatalog.org/score/PGS003378/](http://www.pgscatalog.org/score/PGS003378/) and [www.pgscatalog.org/score/PGS003379/](http://www.pgscatalog.org/score/PGS003379/).

- Genome-wide summary statistics for the PSA multi-ancestry meta-analysis and ancestry-stratified summary statistics for the development of the genome-wide PSA polygenic score are available from: [10.5281/zenodo.7460135](https://doi.org/10.5281/zenodo.7460135).

## Human research participants

Policy information about [studies involving human research participants and Sex and Gender in Research](#).

### Reporting on sex and gender

Because PSA levels are only relevant to individuals with a prostate, all analyses in the study were restricted to individuals of male biological sex (as determined by analyses of genotype data using the Plink software). We additionally restricted to individuals with self-reported male gender identity in an effort to reduce variability in PSA levels attributable to discordance between biological sex and gender identity.

### Population characteristics

- The UK Biobank is a population-based prospective cohort of 502,611 individuals from the United Kingdom, ages 40 to 69 at recruitment between 2006 and 2010. Age at PSA measurement ranged between 21.3 and 78.2 years (mean: 65.5 years). Median PSA across all available values was 2.35 ng/mL (mean: 4.24 ng/mL).

- The Resource for Genetic Epidemiology Research on Aging (GERA) Cohort consists of 100,000 adults who are members of the Kaiser Permanente Medical Care Plan, Northern California Region (KPNC), and participants in its Research Program on Genes, Environment and Health (RPGEH). This analysis used data from men aged between 20 and 90 years (mean: 64.9 years). Median PSA across all available values was 1.4 ng/mL (mean: 4.82 ng/mL).

- The Prostate, Lung, Colorectal and Ovarian (PLCO) Cancer Screening Trial is a completed randomized trial that enrolled approximately 155,000 participants aged 55 to 74 between November 1993 and July 2001. These analyses were limited to men (mean age: 62.5 years) with a baseline PSA measurement who were randomized to the screening arm of the trial (N=29,524) with mean PSA of 1.30 ng/mL (standard deviation = 0.88).

- In BioVU, the mean age at PSA measurement was 56.9 years and the median PSA level across all available measurements was 1.00 ng/mL.

- PCPT randomly assigned 18,880 men aged 55 years or older who had a normal DRE and PSA level  $\leq 3$  ng/mL to either finasteride or placebo. Analyses in this manuscript are based on 5883 participants aged between 55 and 85 years (mean: 63.2 years). Median PSA across all available values was 1.1 ng/mL (mean: 1.21 ng/mL).

- SELECT randomized 34,888 participants aged 50 or older between 2001 and 2004. The minimum enrollment age was 50 years for African American men and 55 years for all other men. Additional eligibility requirements included no prior prostate cancer diagnosis,  $\leq 4$  ng/mL of PSA in serum, and a digital rectal exam not suspicious for cancer. Analyses in this manuscript are based on 25,197 participants aged between 50 and 93 years (mean: 63.0 years). Median PSA across all available values was 1.1 ng/mL (mean: 1.30 ng/mL).

- Details of the participants included in the PSA GWAS from other studies have been previously described: Malmö Diet and Cancer Study (MDCS) [PMID: 22101116]

### Recruitment

Since population-level PSA screening is not currently recommended, all observational PSA data is subject to some level of selection bias arising from patient and healthcare provider preferences, as well as variation in routine clinical care.

The UK Biobank is not representative of the general population across several sociodemographic, physical, lifestyle and health-related characteristics, with evidence of a "healthy volunteer" selection bias, details of which are published elsewhere (Fry et al, *Am J Epidemiol* 2017;186:1026-34. PMID 28641372). Analyses in the presented here are further restricted to a subset of men within the UK Biobank who had linked GP records with available PSA values.

GERA was developed from a mailed survey sent to all adult members of the Kaiser Permanente Medical Care Plan, Northern California Region (KPNC) who had been members for two years or more in 2007. The membership of KPNC is representative of the general population in the 14 county area in which facilities are located, although the membership is underrepresented for the extremes of income at both ends of the spectrum.

The Malmö Diet and Cancer Study (MDCS) is a population-based prospective cohort study that recruited men and women aged between 44 and 74 years old who were living in Malmö, Sweden between 1991 and 1996. These analyses included men from the MDCS who were not diagnosed with prostate cancer as of December 2014 and had available genotyping and baseline PSA measurement.

The Vanderbilt University Medical Center BioVU resource is a synthetic derivative biobank linked to deidentified electronic

health records. Analyses were based on PSA levels that were measured as part of routine clinical care. PLCO enrolled approximately 155,000 participants aged 55 to 74 between November 1993 and July 2001. Men randomized to the screening arm of the trial underwent annual screening with PSA for six years and digital rectal exam (DRE) for four years. These analyses were limited to men randomized to the screening arm of the trial (N= 29,524).

PCPT randomly assigned 18,880 men aged 55 years or older who had a normal DRE and PSA level  $\leq 3$  ng/mL to either finasteride or placebo. Potential biases related to the PCPT design are discussed in detail by Goodman et al. (PMID: 16697846). In SELECT the minimum enrollment age was 50 years for African American men and 55 years for all other men. Additional eligibility requirements included no prior prostate cancer diagnosis,  $\leq 4$  ng/mL of PSA in serum, and a DRE not suspicious for cancer. Analyses presented in this manuscript included PCPT and SELECT participants who were genotyped on the Illumina Infinium Global Screening Array (GSAMD) 24v2-0 array.

#### Ethics oversight

Informed consent was obtained from all study participants. UK Biobank received ethics approval from the Research Ethics Committee (REC reference: 11/NW/0382) in accordance with the UK Biobank Ethics and Governance Framework. The Vanderbilt Institutional Review Board approved the BioVU study. We used previously published PSA GWAS results from the GERA cohort by Hoffmann et al. (PMID: 28139693). The original study was approved by the Kaiser Permanente Northern California Institutional Review Board and the University of California San Francisco Human Research Protection Program Committee on Human Research. The Malmö Diet and Cancer Study (MDCS) was approved by the local ethics committee. The PLCO study was approved by the institutional review board at each participating centre and the National Cancer Institute. The informed consent document signed by the PLCO study participants allows use of these data by investigators for discovery and hypothesis generation in the investigation of the genetic contributions to cancer and other adult diseases. Our study includes publicly posted genomic summary results from PLCO Atlas. No IRB review is required for PLCO summary data use.

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

#### Sample size

Multi-ancestry genome-wide meta-analysis of PSA levels included a total of 95,768 men. This was the final sample size after all relevant exclusions (described below). Results of this analysis were used to develop a PSA genetic score that was validated in individuals that passed quality control in the Prostate Cancer Prevention Trial (n=5883) and the Selenium and Vitamin E Cancer Prevention Trial (n=22,173). These clinical trials have concluded and were converted to observational studies.

#### Data exclusions

Genome-wide association study (GWAS) of PSA levels excluded participants who were ever diagnosed with prostate cancer and individuals with PSA values  $>10$  or PSA=0 (as an indicator of prostate resection). Analyses in the Prostate Lung, Colorectal and Ovarian (PLCO) were limited to men with a baseline PSA measurement who were randomized to the screening arm of the trial. In the UK Biobank we removed participants who withdrew consent at a later date and no longer wish their data to be included. Additional exclusions focused on ensuring that only high-quality genetic data were retained for downstream analyses. Detailed descriptions of the quality control procedures performed by each contributing study are described in the Methods. Briefly, iterative 80% and 95% sample- and variant-level call rate filters were applied to remove poorly genotyped or contaminated samples and variants. Heterozygosity outliers within each ancestral population were detected using absolute values from PLINK method-of-moments F coefficients. Samples with values more than five standard deviations from the population mean were excluded. We also excluded individuals with discordant self-reported and genetically inferred sex based on X chromosome method-of-moments F coefficient from PLINK using 0.5 as the threshold (F coefficients are close to 0.0 for males and 1.0 for females). KING version 2.0 (<http://people.virginia.edu/~wc9c/KING/>) was used to estimate relatedness among the samples based on a subset of genotyped autosomal variants with minor allele frequency (MAF)  $\geq 0.01$  and genotype call rate  $\geq 97\%$ . We excluded one individual from each pair of first-degree relatives. To further minimize potential population stratification in the UK Biobank, we excluded individuals for whom either of the first two genetic ancestry principal components (PC's) were  $>5$  standard deviations away from the mean of the population. GWAS analyses were limited to variants with MAF  $>0.005$  and imputation quality INFO  $>0.30$  in each of the contributing studies. We excluded variants that were out of Hardy-Weinberg equilibrium in cancer-free individuals (p-value  $<1E-05$  in UKB and p-value  $<1E-06$  in some studies).

#### Replication

The goal of the present analysis is to establish the predictive performance the PSA genetic score and the potential clinical utility of using this score to correct PSA measurements. The PSA genetic score was developed from a GWAS of PSA levels in 95,768 men and was subsequently validated in two independent studies that were not part of the GWAS: the Prostate Cancer Prevention Trial (PCPT) and the Selenium and Vitamin E Cancer Prevention Trial (SELECT). We provide the genetic variants and corresponding weights (effect sizes) necessary to construct the PSA genetic score and perform genetic adjustment of PSA levels to facilitate future replication of this work.

#### Randomization

Observational studies (GERA cohort, Malmö Diet & Cancer Study) and biobanks (UK Biobank, BioVU) that contributed data to the PSA GWAS did not have a randomization or intervention component. Analyses in the PLCO were limited to baseline PSA values in the screening arm of the trial. GWAS of PSA levels adjusted for the following minimum set of covariates: age at PSA measurement, the first 10 genetic ancestry principal components, and genotyping array or imputation batch (where applicable). Association analyses of polygenic scores (PGS) and genetically adjusted PSA in relation to prostate cancer incidence that were performed in the Prostate Cancer Prevention Trial (PCPT) and the

Selenium and Vitamin E Cancer Prevention Trial (SELECT) included randomization arm as a covariate in addition to age and genetic ancestry principal components.

## Blinding

All data for performing genome-wide association analyses and developing polygenic scores were de-identified.

The researchers who carried out the analyses for this manuscript had no influence on how the genotyping, PSA measurement, or assessment of cancer status was performed in any of the contributing studies.

Blinding is not relevant for our study because we used observational and EHR/biobank data (GERA, Malmo Diet & Cancer Study, UK Biobank, BioVU). The clinical trials that contributed data to our study have all been completed and converted to observational cohorts: Prostate, Lung, Colorectal and Ovarian (PLCO) Cancer Screening Trial, Prostate Cancer Prevention Trial (PCPT), and the Selenium and Vitamin E Cancer Prevention Trial (SELECT). For the purpose of our manuscript, blinding in these trials is not relevant since our analyses focused on baseline/pre-randomization PSA data for GWAS (in PLCO) and polygenic score validation (in PCPT and SELECT). For associations with prostate cancer in this manuscript, both PCPT and SELECT were analyzed as observational case-control studies.

# Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

## Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

## Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging