

Combined Utility of 25 Disease and Risk Factor Polygenic Risk Scores for Stratifying Risk of All-Cause Mortality

Allison Meisner,¹ Prosenjit Kundu,¹ Yan Dora Zhang,^{1,2} Lauren V. Lan,¹ Sungwon Kim,¹ Disha Ghandwani,^{1,3} Parichoy Pal Choudhury,⁴ Sonja I. Berndt,⁴ Neal D. Freedman,⁴ Montserrat Garcia-Closas,⁴ and Nilanjan Chatterjee^{1,5,*}

While genome-wide association studies have identified susceptibility variants for numerous traits, their combined utility for predicting broad measures of health, such as mortality, remains poorly understood. We used data from the UK Biobank to combine polygenic risk scores (PRS) for 13 diseases and 12 mortality risk factors into sex-specific composite PRS (cPRS). These cPRS were moderately associated with all-cause mortality in independent data within the UK Biobank: the estimated hazard ratios per standard deviation were 1.10 (95% confidence interval: 1.05, 1.16) and 1.15 (1.10, 1.19) for women and men, respectively. Differences in life expectancy between the top and bottom 5% of the cPRS were estimated to be 4.79 (1.76, 7.81) years and 6.75 (4.16, 9.35) years for women and men, respectively. These associations were substantially attenuated after adjusting for non-genetic mortality risk factors measured at study entry (i.e., middle age for most participants). The cPRS may be useful in counseling younger individuals at higher genetic risk of mortality on modification of non-genetic factors.

Introduction

Genome-wide association studies (GWASs) with increasingly large sample sizes have led to the discovery of thousands of genetic variants associated with individual traits, including complex diseases and risk factors for disease.¹ Analyses of polygenicity of a variety of traits^{2,3} have further indicated that many individual traits are likely to be associated with thousands to tens of thousands of genetic variants, each with very small effect. Thus, much attention has been paid to the utility of polygenic risk scores (PRS), which represent the genetic burden of a given trait, for developing strategies for risk-based intervention through lifestyle modification,^{4–8} screening,^{5,7–12} and medication.^{5,7,13,14} A PRS for a given trait is typically defined as a weighted sum of a set of germline single-nucleotide polymorphisms (SNPs), where the weight for each SNP corresponds to an estimate of the strength of association between the SNP and the trait.⁷ Recent studies indicate that while PRS tend to have modest predictive capacity overall, they have the potential to offer substantial stratification of a population into distinct levels of risk for some common diseases such as coronary artery disease (CAD) and breast cancer.^{4,15}

There is ongoing debate regarding the utility of PRS in clinical practice.^{16–18} PRS can be more robust and cost-efficient tools for risk stratification than other biomarkers and risk factors. In particular, PRS do not change over time and thus need to be measured only once. Additionally, the risk associated with PRS for different traits appears in many cases to be fairly consistent over an indi-

vidual's life course^{15,19} and time-varying lifestyle and clinical factors tend to act in a multiplicative way on baseline genetic risk.^{4,6,20,21} Further, if genome-wide genotype and/or sequencing data are available on an individual, the same data can be used to evaluate the PRS for a large number of traits simultaneously. Thus, beyond the use of PRS for prevention of specific diseases, it is important to evaluate their utility for broad health outcomes, particularly if PRS are to be utilized in routine health care.

The broad health impact of public health or clinical interventions is often measured in terms of their impact on all-cause mortality or lifespan.^{22–25} While a small number of genetic variants associated with lifespan have been identified,^{26–28} no study to date has systematically evaluated the ability of emerging PRS for life-threatening diseases and mortality risk factors to predict mortality. We used data from the UK Biobank, a large prospective cohort study, to assess the combined utility of PRS associated with 13 common diseases and 12 established risk factors for mortality. We used training data to combine the trait-specific PRS into sex-specific composite PRS (cPRS) that are predictive of all-cause mortality. We then evaluated the association of these cPRS with all-cause mortality and their ability to stratify mortality risk in independent test data. We also assessed the degree to which mortality risk associated with the cPRS was accounted for by mortality risk factors measured at the time of entry into the study, i.e., middle age for most participants. Finally, we examined the potential clinical use of the cPRS, namely, counseling younger individuals at higher genetic risk of mortality on

¹Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD 21205, USA; ²Department of Statistics, University of Hong Kong, 999077, Hong Kong; ³Indian Statistical Institute, Kolkata, West Bengal 700108, India; ⁴Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda, MD 20892, USA; ⁵Department of Oncology, Johns Hopkins University School of Medicine, Baltimore, MD 21205, USA

*Correspondence: nilanjan@jhu.edu

<https://doi.org/10.1016/j.ajhg.2020.07.002>

© 2020 American Society of Human Genetics.



modification of non-genetic risk factors such as body mass index (BMI) and smoking status.

Subjects and Methods

Causes of Death and Mortality Risk Factors

We used the Centers for Disease Control (CDC) Wide-ranging ONline Data for Epidemiologic Research (WONDER) database to identify the top causes of death (organized by the International Classification of Diseases [ICD]-10 113 Causes List) in terms of the number of deaths among non-Hispanic whites in the United States over the age of 40 in 2017, separately for men and women (see [Web Resources](#)). We then determined the top 10 causes of death with some genetic basis, i.e., causes for which there is evidence of an association between one or more genetic variants and disease risk ([Table S1](#)). These causes accounted for 70.3% and 71.8% of deaths among women and men, respectively, in the CDC data.

Several of these causes were very general categories of disease (e.g., “diseases of heart”), making it difficult to identify relevant trait-specific GWASs. Thus, we identified the specific cause within these categories associated with the highest number of deaths (with the exception of “malignant neoplasms;” here, we identified the top four cancers for each sex in terms of the number of deaths). The final list of diseases was: CAD, COPD, Alzheimer disease, stroke, type 2 diabetes, CKD, hypertension, alcoholic liver cirrhosis, Parkinson disease, pancreatic cancer, colorectal cancer, lung cancer, breast cancer (women only), and prostate cancer (men only) ([Table S1](#)). These causes of death captured 44.4% and 44.9% of deaths among women and men, respectively, in the CDC data. The difference between these figures and those cited above (70.3% and 71.8% for women and men, respectively) are driven largely by deaths from non-CAD diseases of the heart and deaths from malignant neoplasms not included in our list of cancers. As our analysis involves UK Biobank data, we also used Office of National Statistics mortality data (see [Web Resources](#)) to determine the top causes of death in the UK; these were nearly identical to those identified using the CDC data ([Table S1](#)).

Based on government statistics from the UK,²⁹ we further identified major mortality risk factors that are known to have some genetic component.^{30,31} We included smoking status, alcohol consumption, SBP, BMI, total cholesterol, fasting plasma glucose, and eGFR. Beyond the risk factors highlighted by the UK government statistics, we included LDL cholesterol, HDL cholesterol, triglycerides, DBP, and sleep duration. In particular, sleep duration was included on the basis of several studies showing clear links between sleep duration and all-cause mortality.^{32–34}

Extraction of SNP Information from the GWAS Catalog and Publicly Available GWASs

To generate a PRS for each disease included in the top causes of death, we used results published in the NHGRI-EBI GWAS Catalog³⁵ to identify SNPs associated with the disease. We downloaded the GWAS Catalog results on March 15, 2019, and selected autosomal genome-wide significant SNPs (p value $\leq 5 \times 10^{-8}$). For each disease, we identified one or more search terms based on the trait names used by the GWAS Catalog, and we selected the SNPs corresponding to these search terms. We then checked several fields of the GWAS Catalog, such as the source of the data, the study title, and the description of the trait studied, to

ensure that we retained relevant SNPs; in particular, we sought to include results from analyses of Europeans (or multi-ethnic populations including Europeans) and to exclude studies of pleiotropic or composite outcomes, studies not of disease susceptibility, studies of children or pregnant women, studies of a secondary condition in individuals with a primary condition (e.g., myocardial infarction in individuals with coronary heart disease), studies of haplotypes or multi-SNP analyses, and studies of subpopulations (e.g., carriers of a specific genetic mutation; the only exceptions to this were studies of cirrhosis among alcohol drinkers and studies of COPD among smokers) or SNP-environment interactions. Importantly, these exclusions mean we included only GWASs of disease status, rather than GWASs of particular outcomes among individuals with a given disease, e.g., disease-associated mortality. In the resulting list of SNPs, there were several cases where the same SNP appeared multiple times for the same disease trait. In these situations, we kept the result from the largest study (in terms of the number of cases). The same SNP may appear for multiple traits.

For our analysis, it was important to extract the effect allele, effect size, and effect allele frequency for each SNP. The effect allele and effect size were used to construct the PRS in the UK Biobank, and the effect allele and effect allele frequency were used to check whether the SNP in the UK Biobank was the same as the SNP reported on the GWAS Catalog. For many SNPs on the list we created, some or all of this information was missing in the GWAS Catalog. We sought to fill in this information by consulting the original paper and its supplemental materials, as well as the Ensembl database.³⁶ In situations where we were not able to discern the effect allele, the effect allele frequency, or the effect size of a particular SNP, the SNP was removed from our list.

We applied the same approach for identifying SNPs for each cause of death except for stroke. This is because there are several types of stroke and different studies included in the GWAS Catalog employed definitions of stroke with varying specificity. Thus, we used a recently published stroke PRS³⁷ instead of using the results available from the GWAS Catalog.

Our approach to identifying SNPs for inclusion in the mortality risk factor PRS differed from the approach described above. In particular, we found that the risk factor phenotypes were typically defined and/or analyzed differently across studies. For instance, smoking behavior could be defined as ever-use of cigarettes (never versus former/current) or more granularly, incorporating cigarettes per day and duration among ever smokers. As another example, body mass index could be analyzed as a raw measurement, or it could first be rank-transformed. In light of these complications, instead of using the results included in the GWAS Catalog, we used the results from the most recent, largest trait-specific GWAS for which summary data were available (see Neale Lab in [Web Resources](#)).^{38–43} As above, we selected autosomal genome-wide significant SNPs ($p \leq 5 \times 10^{-8}$) and removed SNPs for which the effect allele, effect size, or effect allele frequency were unavailable. In addition, as variant identifiers (RS IDs) were the primary way of querying the UK Biobank genotype data (described below), SNPs without RS IDs were removed (this was not an issue for the GWAS Catalog results).

UK Biobank: Disease and Mortality Data

The UK Biobank is a large cohort study of more than 500,000 individuals in the UK.⁴⁴ The study enrolled individuals aged 40–69 years between 2006 and 2010 and has followed them since

enrollment. A vast array of information has been collected from these individuals, including genotype data, anthropometric measurements, and information on lifestyle factors and personal and family history of disease. Additionally, data from national death and cancer registries are linked to the UK Biobank data.

We retrieved data on mortality, incident, and prevalent disease for the top causes of death, and mortality risk factor measurements at baseline. The death registry data were available through November 30, 2016, for the centers in Scotland and January 31, 2018, for the centers in England and Wales. We determined whether an individual died of a particular disease by considering the ICD-10 code listed as the primary cause of death (see [Table S1](#) for the codes used). We used several sources of data to identify incident and prevalent cases of disease for the top causes of death. In particular, we used cancer registry data (available through October 31, 2015, in Scotland and March 31, 2016, in England and Wales) to determine whether participants had or experienced the cancers in our list of diseases before (prevalent case) or after (incident case) study baseline on the basis of ICD-9 and ICD-10 codes ([Table S2](#)). For the non-cancer diseases, we used questionnaire/interview data, hospital episode data (available through March 31, 2017, in England, October 31, 2016, in Scotland, and February 29, 2016, in Wales), and death registry data to identify prevalent and incident cases of disease ([Table S2](#)). The exception to this was incident and prevalent diabetes, which were defined based on the algorithm presented in Eastwood et al.⁴⁵ For SBP and DBP at baseline, two measurements were made for each; when both of these were non-missing, the average was used. Self-reported intake of different forms of alcohol was converted into grams of alcohol per day ([Table S3](#)).

In all analyses, unless otherwise specified, we adjusted for the first ten genetic principal components, which were provided by the UK Biobank, in order to account for population stratification. In addition, all survival models accounted for left truncation by starting the follow-up interval at study entry. Throughout, we restricted our attention to unrelated participants (third degree relatives or closer were removed) of white British ancestry, in order to minimize the influence of population stratification and avoid issues related to clustering of individuals in families. We further removed individuals who had withdrawn their consent to participate. Unrelated participants were identified as those who were used by the UK Biobank to compute the principal components and ancestry was determined by the UK Biobank based on self-report and principal component analysis. The UK Biobank was approved by the North West Multi-centre Research Ethics Committee. This research was conducted using the UK Biobank Resource under Application Number 17712.

Evaluating PRS in the UK Biobank

Imputed genotype data (in the form of allele dosage, i.e., between 0 and 2) for the SNPs identified above were extracted from the UK Biobank, matching on RS ID if possible and on chromosome and position otherwise. Non-bi-allelic SNPs and ambiguous palindromic SNPs (A/T or C/G SNPs with allele frequencies between 0.4 and 0.6) were removed. To ensure the SNPs from the UK Biobank were the same as those on our curated list of trait-associated SNPs, the alleles and allele frequencies were compared (allowing for the possibility of strand flips). SNPs that did not match the UK Biobank data, i.e., SNPs for which the reported allele frequency and the allele frequency in the UK Biobank differed by more than 0.15, were removed. Finally, SNPs in linkage disequilibrium (LD)

were removed via LD clumping, implemented using PLINK with an r^2 cutoff of 0.1 and based on the reported p values (from the GWAS Catalog or the publicly available summary statistics) and the 1000 Genomes European reference panel.^{46,47} This was done separately for each disease and risk factor, yielding a list of independent SNPs for each trait. The one exception was stroke: the SNP list was not pruned because the estimated association coefficients provided were based on a joint SNP model. The number of SNPs included in each PRS varied widely, between two SNPs for cirrhosis and 1,458 for BMI ([Table S4](#)). In total, our analysis included 3,941 unique SNPs.

Next, a PRS for each trait was constructed for each participant by weighting the SNP dosage by the reported log odds ratio (for binary traits) or linear regression coefficient (for continuous traits):

$$PRS_{i,j} = \sum_{k=1}^{m_j} g_{i,k} \beta_{k,j},$$

where $PRS_{i,j}$ is the PRS value for the i^{th} individual and the j^{th} trait, m_j is the number of SNPs included in the PRS for the j^{th} trait, $g_{i,k}$ is the genotype dosage for the i^{th} individual and the k^{th} SNP, and $\beta_{k,j}$ is the log odds ratio or linear regression coefficient for the k^{th} SNP and the j^{th} trait.

Statistical Analysis

All analyses were sex-specific and the PRS were standardized to have unit variance. We first evaluated the association between each derived PRS and the corresponding trait (i.e., prevalent disease and incident disease for the disease trait, and measurement at baseline for the mortality risk factors). For the disease traits, we evaluated the association with incident and prevalent disease status separately. To evaluate the relationship between each disease PRS and prevalent disease, we fit a logistic regression model for each disease. We used Poisson models with robust variance estimation⁴⁸ to evaluate the association between each disease PRS and incident disease among individuals without prevalent disease. For the mortality risk factors, we used linear regression with robust variance estimation to model the relationship between each mortality risk factor PRS and the risk factor measurement at baseline. The one exception was smoking status; since the smoking status PRS was developed based on a GWAS of ever-use of cigarettes, we defined the smoking status risk factor as ever-use of cigarettes. As this is a binary variable, we used logistic regression to model the relationship between the smoking status PRS and ever-use of cigarettes. Since eGFR was not directly available in the UK Biobank, we calculated eGFR at baseline using the Modification of Diet in Renal Disease (MDRD) Study Equation;⁴⁹ this mirrors the definition of eGFR used in the GWAS upon which our eGFR PRS was based.⁴² All models included adjustment for age at entry, in addition to the first ten principal components.

We also investigated cause-specific mortality for the diseases included in our top causes of death. We used Cox proportional hazards models to study the relationship between each disease PRS and age at death from that disease. Deaths from other causes were treated as censoring events. We performed these analyses in the full cohort and also among individuals with and without the disease corresponding to the cause of death being modeled at baseline. We also evaluated the relationship between each mortality risk factor PRS and mortality due to each of the causes of death. For all of the analyses related to cause-specific mortality, when there were not enough deaths to yield stable estimates, estimates are not provided.

Table 1. Descriptive Statistics

	Full Cohort		Training Data		Test Data	
	Women	Men	Women	Men	Women	Men
Sample size	181,027	156,111	120,719	104,037	60,308	52,074
Age at study entry (years; mean [SD])	57.2 (7.9)	57.6 (8.1)	57.2 (7.9)	57.6 (8.1)	57.2 (7.9)	57.6 (8.1)
Follow-up (years; mean [SD])	8.8 (1.1)	8.7 (1.3)	8.8 (1.1)	8.7 (1.3)	8.8 (1.0)	8.7 (1.3)
Number of deaths	5,250	8,360	3,530	5,576	1,720	2,784

SD, standard deviation. Descriptive statistics for the full cohort used for the analysis (after removing individuals who were related, were not of British ancestry, or had withdrawn their consent to participate), the training data (2/3 of the full cohort), and the test data (1/3 of the full cohort).

Our main analysis involved studying the joint relationship between the 25 PRS and all-cause mortality. First, we split the data into training (2/3) and test (1/3) sets. Then, in the training data, all PRS (with the exception of prostate cancer and breast cancer for the female- and male-specific models, respectively) were included in Cox proportional hazards models of age at death:

$$\lambda(t|PRS_1, \dots, PRS_{25}, \mathbf{Z}) = \lambda_0(t) \exp(\theta_1 PRS_1 + \dots + \theta_{25} PRS_{25} + \beta^T \mathbf{Z}).$$

In this formula, $\lambda(t|PRS_1, \dots, PRS_{25}, \mathbf{Z})$ denotes the hazard at age t given $(PRS_1, \dots, PRS_{25}, \mathbf{Z})$, $\lambda_0(t)$ denotes the baseline hazard at age t , and \mathbf{Z} is a vector of the first ten principal components. Each model yielded a weighted combination of the individual PRS where the weights were the estimated log hazard ratios (HRs) from the Cox model, $\hat{\theta}_1 PRS_1 + \dots + \hat{\theta}_{25} PRS_{25}$; we refer to these sex-specific weighted combinations as the “composite PRS” (cPRS). These cPRS were then applied to the test data. In particular, we used a Cox model to evaluate the HR for all-cause mortality per standard deviation of the cPRS. In addition, we estimated the HR comparing individuals in the top 5% of the cPRS distribution to those in the middle 20% and the HR comparing individuals in the bottom 5% to those in the middle 20% in the test data. This was based on quantiles estimated in the training data. To aid in the interpretation of these results, the estimated HRs were converted into approximate years of life difference, as done in other studies of survival.^{26,31} In addition, we used Harrell’s C-index to quantify the discriminatory ability of the cPRS;⁵⁰ note that this evaluation did not adjust for principal components.

We undertook a series of additional analyses. First, we evaluated the association between the cPRS and all-cause mortality in the “healthy” subset of the test data, that is, the test set after removing individuals with any of the diseases included as a top cause of death at baseline (i.e., prevalent cases). We also re-evaluated the association between the cPRS and all-cause mortality in the test data, adjusting for the mortality risk factors measured at baseline (that is, BMI, smoking status, alcohol consumption, SBP, DBP, eGFR, total cholesterol, LDL cholesterol, HDL cholesterol, triglycerides, blood glucose, and sleep duration), removing individuals in the test data that were missing any of these measurements. All risk factors were included as continuous variables, with the exception of smoking status, which was included as a binary variable (ever versus never use).

Finally, we evaluated the relationship between two major modifiable risk factors, BMI and smoking status, and absolute risk of mortality for individuals at different levels of polygenic risk. We estimated the mortality risk for obese individuals (BMI > 30 kg/m²) and normal weight individuals (BMI of 18.5–25 kg/m²) based on Cox proportional hazards models with quintiles of the

cPRS and BMI categories (≤ 18.5 kg/m², (18.5–25 kg/m²), (25–30 kg/m²), >30 kg/m²), both modeled as categorical variables, fit in the test data. Estimates of risk for never smokers and ever smokers are based on Cox proportional hazards models with quintiles of the cPRS, modeled as a categorical variable, and an indicator of ever-use of cigarettes, fit in the test data. These models did not include adjustment for principal components.

All analyses were conducted using R, including the rms, survival, ggplot2,⁵¹ and sandwich^{52,53} packages (see [Web Resources](#)). We report 95% confidence intervals throughout.

Results

UK Biobank: Disease, Mortality, and Genotype Data

After removing individuals who were related, were not of British ancestry, or had withdrawn their consent to participate, our dataset included 337,138 participants, including 181,027 women and 156,111 men (Tables 1 and S5). There were 13,610 deaths (4.0%) with 5,250 among women (2.9%) and 8,360 among men (5.4%). The diseases included in the top causes of death accounted for 45.9% of the deaths in women and 45.5% of the deaths in men in the UK Biobank. Notably, very few deaths in the UK Biobank were attributed to type 2 diabetes, which appears to be due to many more deaths in the UK Biobank having type 2 diabetes listed as a secondary cause of death as opposed to the primary cause.

Constructing and Evaluating the Trait-Specific PRS in the UK Biobank

As anticipated, the trait-specific PRS tended to be moderately to strongly associated with the corresponding disease or risk factor (Figure S1 and Table S6). The strongest associations for the disease traits (odds ratios or relative risks of at least 1.5 per standard deviation [SD]) were observed for Alzheimer disease (incident disease only), type 2 diabetes, breast cancer in women, prevalent CAD in men, cirrhosis in men, and prostate cancer in men.

We observed that the PRS for each disease was generally at least moderately associated with death from that disease (Figure 1), with the association being strongest for Alzheimer disease (HR per SD: 1.86 [95% confidence interval: 1.42, 2.42] in women; 2.01 [1.52, 2.65] in men), CAD (1.51 [1.34, 1.69] in women; 1.48 [1.40, 1.57] in men),

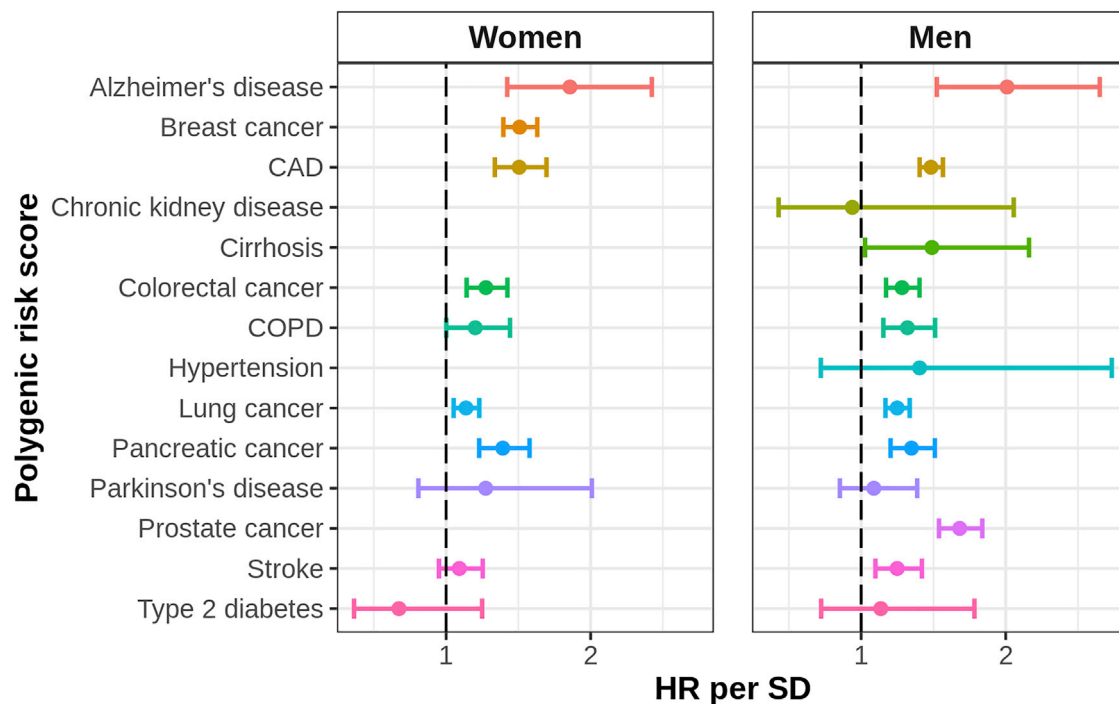


Figure 1. Association of Each Disease PRS with Cause-Specific Mortality in the Full Cohort

For each disease, we evaluated the association between the disease PRS and mortality from the disease based on sex-specific Cox proportional hazards models of age at death. Deaths from other causes were treated as censoring events. Some causes did not have enough deaths to yield stable estimates, i.e., <6 deaths; in these cases, estimates are not provided. Each PRS was standardized to have unit variance so the estimates correspond to the HR per SD of the PRS. The horizontal lines indicate 95% confidence intervals. CAD, coronary artery disease; COPD, chronic obstructive pulmonary disease; HR, hazard ratio; SD, standard deviation; PRS, polygenic risk score.

breast cancer in women (1.51 [1.40, 1.63]), prostate cancer in men (1.68 [1.54, 1.84]), and cirrhosis in men (1.49 [1.03, 2.16]). In general, the PRS were stronger predictors of cause-specific mortality among individuals without prevalent disease than they were among individuals with prevalent disease (Figure S2); this indicates the PRS were typically more strongly associated with disease onset than with prognosis.

We found that the PRS for BMI was at least moderately associated with mortality related to CAD (primarily in men), COPD (among women), hypertension (among men), lung cancer (among women), pancreatic cancer (among women), Parkinson disease (among women), and stroke (among women) (Figures S3 and S4). The PRS for smoking was weakly associated with mortality due to CAD (among men) and moderately associated with mortality due to COPD (primarily in men) and lung cancer. The PRS for LDL cholesterol was strongly associated with mortality related to Alzheimer disease (among men) and COPD (among women) and moderately associated with mortality due to CAD (primarily in men). The PRS for total cholesterol was strongly positively associated with mortality due to Alzheimer disease (primarily in men) and COPD (among women), moderately positively associated with mortality related to CAD (among men), and moderately negatively associated with mortality due to pancreatic cancer (among men). The PRS for triglycerides was strongly negatively associated with mortality from stroke among

men. The PRS for alcohol consumption was moderately positively associated with mortality due to CAD, primarily among men.

We found that several PRS were modestly associated with all-cause mortality, with some differences between men and women (Figure 2). The PRS for BMI was modestly associated with risk of all-cause mortality for both women (HR per SD: 1.07 [1.04, 1.10]) and men (1.08 [1.05, 1.10]). In addition, the PRS for smoking status, Alzheimer disease, LDL cholesterol, and lung cancer were modestly associated with all-cause mortality in both sexes. The PRS for breast cancer and prostate cancer were modestly associated with all-cause mortality in women and men, respectively. Among men, the PRS for CAD, cirrhosis, DBP, HDL cholesterol, SBP, stroke, total cholesterol, triglycerides, type 2 diabetes, and alcohol consumption were modestly associated with all-cause mortality; notably, the PRS for HDL cholesterol and triglycerides were both negatively associated with all-cause mortality. In general, the estimated associations tended to be stronger in men than in women.

Constructing and Evaluating the Composite PRS in the UK Biobank

The training data used to construct the cPRS included 224,756 participants, among them 120,719 women and 104,037 men (Table 1). There were 9,106 deaths in the training data with 3,530 in women and 5,576 in men. Correspondingly, the test data used to evaluate the cPRS

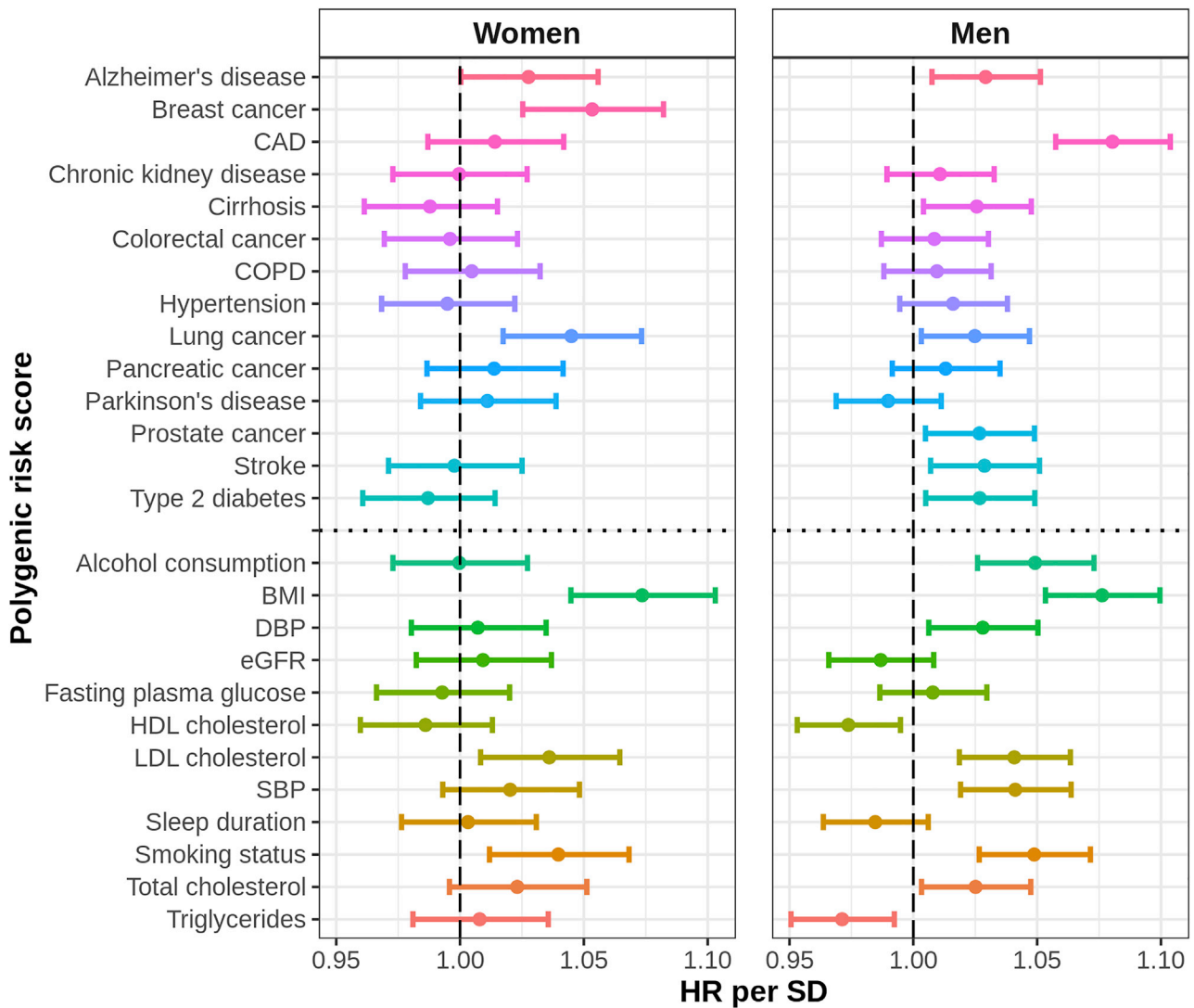


Figure 2. Association of Each Trait-Specific PRS with All-Cause Mortality in the Full Cohort

We evaluated the association between each PRS and all-cause mortality based on sex-specific Cox proportional hazards models of age at death in the full cohort. Each Cox model included one PRS. Each PRS was standardized to have unit variance so the estimates correspond to the HR per SD of the PRS. The horizontal lines indicate 95% confidence intervals. BMI, body mass index; CAD, coronary artery disease; COPD, chronic obstructive pulmonary disease; DBP, diastolic blood pressure; eGFR, estimated glomerular filtration rate; HDL, high-density lipoprotein; LDL, low-density lipoprotein; SBP, systolic blood pressure; HR, hazard ratio; SD, standard deviation; PRS, polygenic risk score.

included 112,382 individuals (60,308 women and 52,074 men) and 4,504 deaths (1,720 among women and 2,784 among men).

The cPRS were moderately associated with all-cause mortality in the test data (HR per SD: 1.10 [1.05, 1.16] in women, 1.15 [1.10, 1.19] in men; see [Table 2](#) and [Figure S5](#)). However, the cPRS were able to identify substantial fractions of the population that have meaningfully elevated and reduced mortality risk, particularly among men ([Table 2](#) and [Figure 3](#)). The estimated difference in life expectancy between the top and bottom 5% of the cPRS distribution was 4.79 (1.76, 7.81) years in women and 6.75 (4.16, 9.35) years in men. The overall discriminatory capacity of the cPRS, measured by Harrell's C-index,⁵⁰ was small: 0.525 in women and 0.536 in men. These are

comparable to the values for several strong risk factors for mortality, including BMI (0.532 in women, 0.530 in men), smoking status (0.562 in women, 0.574 in men), and alcohol consumption (0.509 in women, 0.547 in men).

When we evaluated the cPRS in the "healthy" subset of the test data, the estimated associations between the cPRS and all-cause mortality were fairly similar to the results from the main analysis ([Table S7](#)). Separately, when we adjusted for the mortality risk factors measured at baseline, the association between the cPRS and all-cause mortality was markedly attenuated for both sexes ([Table 2](#)). These results indicate that a substantial fraction (40.7% for women and 32.5% for men) of the association between the cPRS and all-cause mortality was accounted for by

Table 2. The Results of the Main Analysis of All-Cause Mortality and the cPRS, with and without Adjustment for Mortality Risk Factors

	Women	Men
Without Adjustment for Mortality Risk Factors		
Population in test data: N (deaths)		
Total population	60,308 (1,720)	52,074 (2,784)
Top 5% of cPRS	3,060 (107)	2,454 (159)
Middle 20% of cPRS	12,005 (342)	10,387 (539)
Bottom 5% of cPRS	3,096 (69)	2,526 (89)
cPRS: HR (95% CI)		
Per SD of cPRS	1.10 (1.05, 1.16)	1.15 (1.10, 1.19)
Top 5% versus middle 20% of cPRS	1.24 (1.00, 1.54)	1.27 (1.07, 1.52)
Bottom 5% versus middle 20% of cPRS	0.77 (0.59, 1.00)	0.65 (0.52, 0.81)
Top 5% versus bottom 5% of cPRS	1.61 (1.19, 2.18)	1.96 (1.52, 2.55)
cPRS: years of life lost (95% CI)		
Per SD of cPRS	0.97 (0.50, 1.44)	1.36 (0.98, 1.73)
Top 5% versus middle 20% of cPRS	2.17 (0.00, 4.34)	2.42 (0.65, 4.19)
Bottom 5% versus middle 20% of cPRS	-2.61 (-5.20, -0.03)	-4.33 (-6.58, -2.09)
Top 5% versus bottom 5% of cPRS	4.79 (1.76, 7.81)	6.75 (4.16, 9.35)
With Adjustment for Mortality Risk Factors		
Population in test data: N (deaths)		
Total population	36,008 (855)	36,283 (1,730)
Top 5% of cPRS	1,799 (51)	1,689 (102)
Middle 20% of cPRS	7,143 (168)	7,240 (329)
Bottom 5% of cPRS	1,907 (37)	1,804 (60)
cPRS: HR (95% CI)		
Per SD of cPRS	1.06 (0.99, 1.13)	1.10 (1.04, 1.15)
Top 5% versus middle 20% of cPRS	1.19 (0.87, 1.63)	1.25 (1.00, 1.56)
Bottom 5% versus middle 20% of cPRS	0.88 (0.62, 1.26)	0.73 (0.55, 0.96)
Top 5% versus bottom 5% of cPRS	1.35 (0.88, 2.07)	1.71 (1.24, 2.36)
cPRS: years of life lost (95% CI)		
Per SD of cPRS	0.58 (-0.11, 1.26)	0.92 (0.43, 1.40)
Top 5% versus middle 20% of cPRS	1.72 (-1.43, 4.86)	2.20 (-0.03, 4.43)
Bottom 5% versus middle 20% of cPRS	-1.27 (-4.85, 2.30)	-3.19 (-5.95, -0.43)
Top 5% versus bottom 5% of cPRS	2.99 (-1.28, 7.26)	5.39 (2.18, 8.60)

BMI, body mass index; CI, confidence interval; cPRS, composite polygenic risk score; DBP, diastolic blood pressure; eGFR, estimated glomerular filtration rate; HDL, high-density lipoprotein; HR, hazard ratio; LDL, low-density lipoprotein; SBP, systolic blood pressure; SD, standard deviation. The cPRS were constructed in the training data and evaluated by fitting sex-specific Cox proportional hazards models of the association between the cPRS and age at death from all causes in the test data. Both the continuous cPRS and categorical cPRS were modeled. The estimated HRs and CIs were converted to estimated years of life lost. The analysis adjusting for mortality risk factors included adjustment for the risk factors measured at baseline (BMI, smoking status, alcohol consumption, SBP, DBP, eGFR, total cholesterol, LDL cholesterol, HDL cholesterol, triglycerides, blood glucose, and sleep duration); individuals missing any of these measurements were excluded.

these risk factors, which are (to varying degrees) heritable traits. After controlling for the measured risk factors, the difference in life expectancy between the top 5% and the bottom 5% of the cPRS distribution was estimated to be 2.99 (-1.28, 7.26) years in women and 5.39 (2.18, 8.60) years in men.

Finally, we evaluated the relationship between BMI and smoking status and absolute risk of mortality for individuals at different levels of polygenic risk (Figure 4). We observe that the estimated 10-year absolute risk of mortality for a 60-year-old woman in the top 20% of the cPRS distribution who is obese is 0.044. This is 38% higher than the

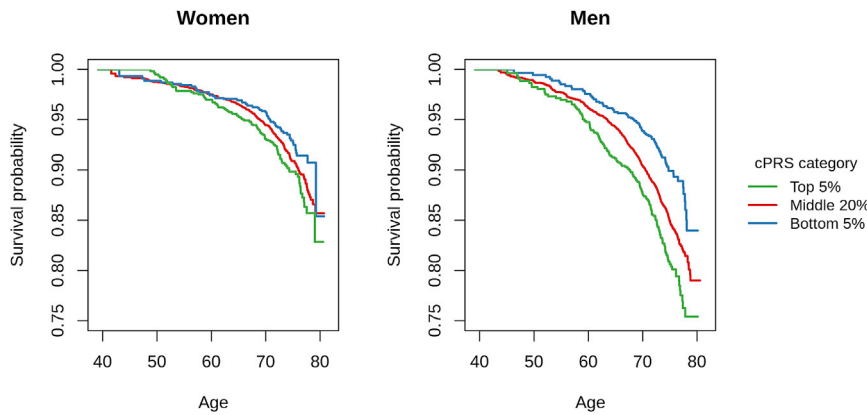


Figure 3. Kaplan-Meier Survival Curves by Quantile of the cPRS

We estimated the sex-specific Kaplan-Meier survival curves for all-cause mortality by quantile of the cPRS in the test data. The Kaplan-Meier curves do not include adjustment for principal components. cPRS, composite polygenic risk score.

directly from the individual SNPs, counting only the number of detrimental or protective alleles across the variants (i.e., without weighting the SNPs by the strength of association). In a combined analysis of men

and women from two studies of northern European populations, the study reported a 10% higher risk of mortality between individuals in the 4th versus 1st quartile of the resulting PRS. In contrast, in the current study, we focus on a limited number of the most important causes of and risk factors for mortality and build cPRS for mortality based on the underlying PRS. Our cPRS, although evaluated in a different population, appears to provide greater mortality risk stratification (HR for 4th versus 1st quartile = 1.29 [1.13, 1.48] in women; 1.38 [1.24, 1.53] in men). These differences may be due to the incorporation of a larger number of SNPs emerging from more recent GWASs as well as the weighting of individual SNPs to account for their association with the individual diseases and risk factors in our analysis.

estimated risk for a woman in the top 20% of the cPRS distribution who is not obese. Similarly, the estimated risk for a 60-year-old woman in the top 20% of the cPRS distribution who is a current or former smoker is 64% higher than for a woman who has never smoked (0.046 versus 0.028). Likewise, for a 60-year-old man, the estimated 10-year risk of mortality is 24% higher if the man is obese as opposed to normal weight (0.087 versus 0.070) and the estimated risk is 81% higher if the man is a current or former smoker relative to a man who has never smoked (0.087 versus 0.048). These differences highlight the potential importance of lifestyle modification even among those at high genetic risk. Furthermore, in most of these examples, the estimated risk for an individual who is in the top 20% of the cPRS distribution but who has a favorable risk factor profile is below the estimated risk for an individual in the middle 20% of the cPRS distribution, i.e., someone at moderate genetic risk (0.032 in women and 0.059 in men).

Discussion

Analyses using a large dataset from the UK Biobank indicate that sex-specific composite PRS (cPRS) for all-cause mortality have fairly modest predictive capacity overall. However, there is evidence that the cPRS could identify substantial fractions of the population with notably elevated and reduced risk of all-cause mortality due to the genetic risk accumulated across many variants. Importantly, our results also show that a substantial proportion of the association between the cPRS and mortality was accounted for by mortality risk factors measured in middle age. These findings suggest that those individuals at high genetic risk of mortality may derive substantial benefit from modification of lifestyle factors; in particular, the cPRS could be useful in counseling individuals at high genetic risk on possible lifestyle choices that are associated with lower mortality risk.

A previous study evaluated the utility of 707 SNPs identified from GWASs of 125 diseases and risk factors for estimating mortality risk.³⁰ This study developed a PRS

Several recent studies^{26,54–57} have investigated the association of individual genetic variants and PRS with parental lifespans due to the increased power of these analyses relative to analyses of lifespan in genotyped individuals. Two large GWASs of parental lifespan, both including data from the UK Biobank, identified a total of only 18 loci,^{26–28} highlighting major challenges in finding individual variants related to lifespan. We constructed a lifespan PRS based on 17 of these variants (one was excluded as it was a palindromic SNP whose direction could not be resolved) and found modest associations with all-cause mortality (HR per SD: 1.02 [0.99, 1.05] in women and 1.04 [1.02, 1.06] in men). We further constructed a new cPRS, which included the 25 disease and risk factor PRS constructed for our analysis as well as the lifespan PRS; the associations of this new cPRS with all-cause mortality were nearly identical to that of the original cPRS (HR per SD of the new cPRS: 1.10 [1.05, 1.15] in women and 1.14 [1.10, 1.19] in men).

An important limitation of previous studies is the lack of adjustment for known mortality risk factors in characterizing the potential utility of PRS for estimating mortality risk. In our analysis, the association between the cPRS and mortality was attenuated by more than 30% after adjusting for the mortality risk factors under study. These results suggest that while genetic variants associated with complex traits in GWASs could provide some mortality

and women from two studies of northern European populations, the study reported a 10% higher risk of mortality between individuals in the 4th versus 1st quartile of the resulting PRS. In contrast, in the current study, we focus on a limited number of the most important causes of and risk factors for mortality and build cPRS for mortality based on the underlying PRS. Our cPRS, although evaluated in a different population, appears to provide greater mortality risk stratification (HR for 4th versus 1st quartile = 1.29 [1.13, 1.48] in women; 1.38 [1.24, 1.53] in men). These differences may be due to the incorporation of a larger number of SNPs emerging from more recent GWASs as well as the weighting of individual SNPs to account for their association with the individual diseases and risk factors in our analysis.

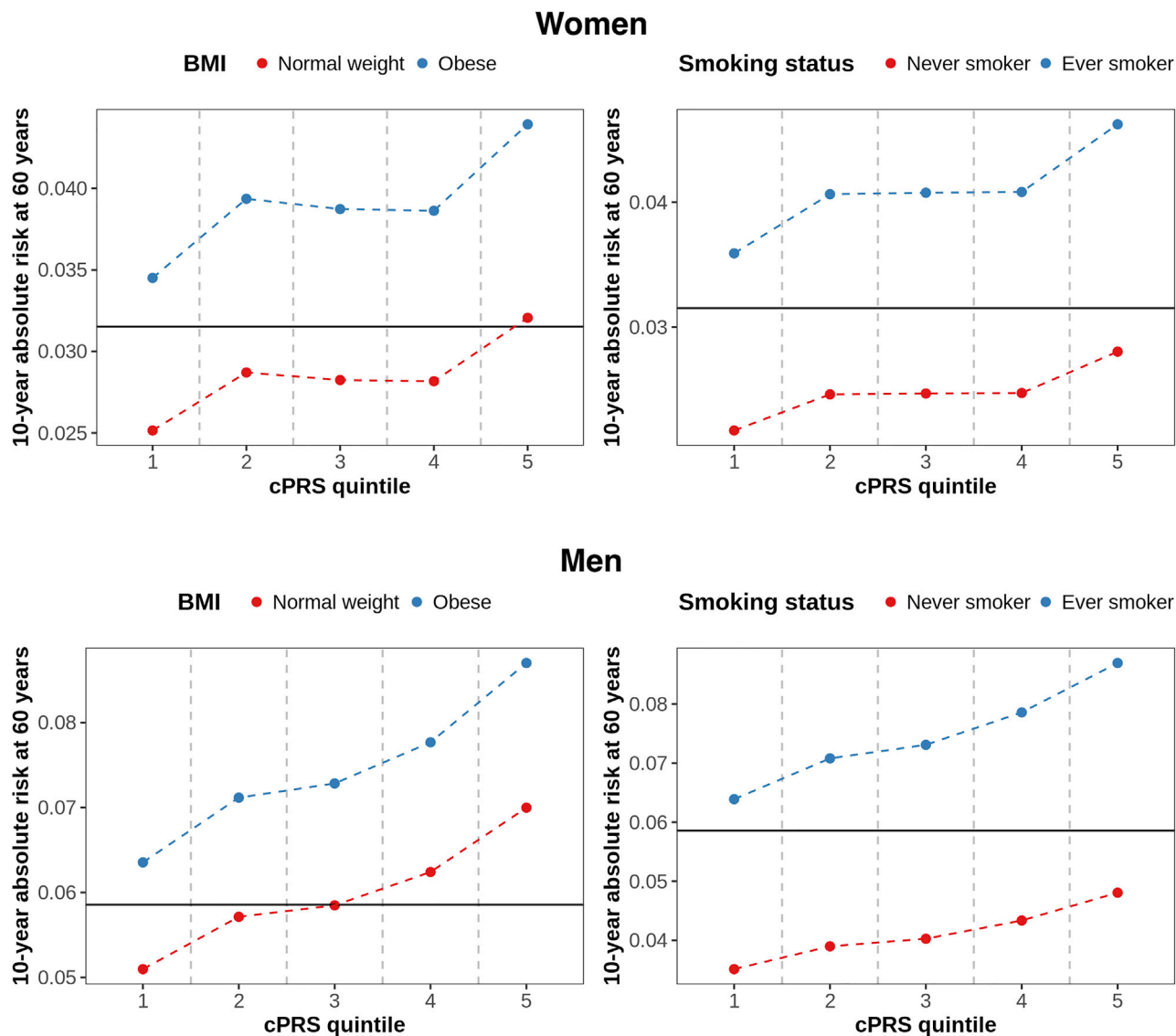


Figure 4. Estimates of Absolute Risk of Mortality in Different Strata of the cPRS within Specific Categories of BMI and by Smoking Status We generated estimates of 10-year absolute risk of all-cause mortality for a 60-year-old in different strata of the cPRS for specific values of two mortality risk factors, BMI and smoking status. The results for women and men are presented in the top and bottom rows, respectively. The horizontal line in each plot corresponds to an estimate of 10-year absolute risk of all-cause mortality for a 60-year-old in the middle quintile of the cPRS, based on sex-specific Cox proportional hazards models with quintiles of the cPRS, modeled as a categorical variable, fit in the test data. BMI, body mass index; cPRS, composite polygenic risk score.

risk stratification early in life, their utility later in life, when other risk factors for mortality can be measured, is diminished.

Most GWASs are case-control studies of disease risk as opposed to prognosis, i.e., aggressiveness and/or progression of the disease leading to death. When we examined the association of the disease PRS with the corresponding cause-specific mortality among individuals with prevalent disease in the UK Biobank (Figure S2), only the PRS for CAD and COPD were (at least moderately) associated; in other words, for most PRS, there was little to no evidence of an association with prognosis or disease survival. Although such analyses may be influenced by selection associated with survivorship and poor health, in general,

there is little evidence of association between disease risk SNPs (and thus disease PRS) and survival following disease onset. While future GWASs focusing on genetic determinants of aggressiveness and disease progression are needed, finding associations may be challenging due to available sample sizes and heterogeneity as a result of various factors such as treatment.

Our analysis of the relationship between the individual PRS and all-cause mortality revealed some important patterns (Figure 2). The strongest positive associations were seen for the PRS for BMI, breast cancer (in women), CAD (in men), smoking status (particularly in men), and alcohol consumption (in men). In addition, weaker associations with all-cause mortality were seen for the PRS for

Alzheimer disease, lung cancer, and LDL cholesterol in both sexes and, among men, associations were seen for the PRS for stroke, cirrhosis, total and HDL cholesterol, prostate cancer, triglycerides, SBP, DBP, and type 2 diabetes. The negative association observed among men for the triglycerides PRS appears to be driven by a strong negative association between the triglycerides PRS and stroke-specific mortality (Figure S4), which is consistent with the “triglycerides paradox” reported by others.^{58–61}

Given that the associations of the CAD PRS with CAD-specific mortality were similar for men and women, the differences in the associations with all-cause mortality may be due to lower rates of CAD in women during the relatively short follow-up period of the UK Biobank. Differential event rates for some diseases for which alcohol consumption is a risk factor (e.g., CAD) could also partially explain the differences observed in the association of the alcohol consumption PRS with all-cause mortality by sex. We note that the sex differences observed in our results more generally are supported by other studies, which have similarly found indications of differences between men and women in the mechanisms governing lifespan and longevity.^{26,27,31,55,56,62,63}

Our results are generally consistent with a recent paper looking at PRS for many clinical risk factors and mortality across the UK Biobank, a Finnish biobank (FinnGen), and Biobank Japan.⁶⁴ In this multi-ethnic study, several modest associations were observed, including for the PRS for SBP, DBP, and BMI (HRs of around 1.03–1.04 per SD in the trans-ethnic meta-analysis). Interestingly, the results from this analysis varied by ethnicity: for instance, within the UK Biobank, the association between the PRS for BMI and mortality reported in Sakaue et al.⁶⁴ was stronger than was observed in the trans-ethnic meta-analysis (HR of approximately 1.07 per SD in the UK Biobank versus 1.04 in the meta-analysis). This highlights the importance of multi-ethnic analyses.

We evaluated the broad utility of PRS in terms of their combined ability to predict mortality. In the future, other broad measures of health outcomes and expenditures, such as disability-adjusted life years (DALYs), should also be considered. The framework we have created for combining individual PRS could be used to create composite PRS for DALYs or other measures. Given that PRS are known to be strongly associated with incidence of many debilitating diseases, one would anticipate such a composite PRS will have greater utility for predicting DALYs than for mortality. However, analysis of DALYs in a cohort study with limited follow-up, like the UK Biobank, is challenging.

Our analysis has several strengths. We used data from the UK Biobank, a large cohort study, to carry out a comprehensive analysis of PRS for complex traits and mortality, both overall and cause-specific. We used a novel approach to derive composite PRS across many diseases and risk factors to evaluate their combined utility for predicting overall mortality. Under the assumption that com-

mon genetic variants identified through recent GWASs influence mortality risk through the outcomes underlying the GWASs, the composite PRS approach provides a more parsimonious and powerful approach to building models for predicting composite outcomes than building models based on individual SNPs. The weights of individual SNPs in a PRS account for the strength and direction of association of each SNP with the corresponding outcome and the weights for the individual PRS in the cPRS reflect (in part) the relative contribution of the individual diseases and risk factors to mortality. Further, we conducted an unbiased evaluation of the performance of the cPRS for predicting mortality by building it in a training dataset and evaluating it in an independent test dataset.

As the UK Biobank participants are volunteers, there is evidence that this cohort differs from the general UK population in important ways, including being less likely to be obese, smoke, or drink alcohol.⁶⁵ Selection bias,⁶⁵ which contributes to such differences, could influence the generalizability of our results.⁶⁶ Additionally, while our cPRS include germline mutations and so could potentially be evaluated at birth, the UK Biobank is comprised of individuals who have survived to at least middle age. Consequently, the results may not be fully generalizable to younger individuals and must be validated in other populations. Furthermore, the analysis of the cPRS with adjustment for the mortality risk factors required excluding observations in the test data with missing values for any of these risk factors. These observations constituted a substantial portion of the test data (40.3% in women, 30.3% in men). However, as the missingness mechanism for at least some risk factors is expected to be not random (e.g., individuals choosing not to answer questions regarding smoking status or alcohol consumption due to the social stigma surrounding these behaviors), imputation is not appropriate. Thus, some caution is warranted in interpreting these results.

As our analysis involved the evaluation of a large number of associations, issues related to multiple comparisons are a potential concern. However, our main analysis of the cPRS was carefully defined *a priori* and performed in independent test data. The other analyses we performed were intended to check the validity of the PRS we developed and to better understand the results of the main analysis of the cPRS. Additionally, we emphasize the strength of association rather than statistical significance in interpreting the results throughout.

Another potential limitation of this analysis was our use of the GWAS Catalog to identify SNPs for inclusion in the disease PRS. As the GWAS Catalog is not an exhaustive listing of SNPs associated with every trait, we may have missed some associated SNPs. Further, our approach to constructing PRS for the mortality risk factors involved a fairly simple approach based on summary statistics. To investigate the degree to which better performance could be achieved by using a more sophisticated approach to constructing PRS, we considered two additional methods

for 14 traits for which summary statistics were available: (1) LD clumping and thresholding (with various LD and p value thresholds)⁶⁷ and (2) LDpred, a Bayesian method that incorporates information on LD structure⁶⁸ (see [Supplemental Material and Methods](#)). We evaluated the predictive capacity of the PRS constructed by the three methods for each trait (i.e., the relationship between the PRS and the corresponding trait) and found that the PRS built using the GWAS Catalog generally performed well. In the case of breast cancer and stroke, however, small but meaningful gains in performance were seen for the LDpred PRS. We repeated the main composite PRS mortality analysis, where the LDpred PRS for breast cancer and stroke were used in place of the corresponding GWAS Catalog PRS to construct a new composite PRS in the training data. We found a stronger association between this new composite PRS and all-cause mortality in the test data among women (HR per SD: 1.14 [95% confidence interval: 1.09, 1.19] versus 1.10 [1.05, 1.16]) and a similar association among men (1.16 [1.11, 1.20] versus 1.15 [1.10, 1.19]).

These results indicate that our approach, which allowed us to apply a uniform procedure for SNP selection to all traits (regardless of whether summary statistics were available), captured most of the known genetic susceptibility for nearly all traits considered. Importantly, these results also suggest that as the ability of trait-specific PRS to predict the corresponding trait continues to improve (driven by increasing GWAS power and novel methods for PRS construction), their utility in stratifying mortality risk is also likely to increase. On the other hand, the ability of trait-specific PRS, regardless of the method used to construct them, to predict all-cause mortality will be limited by the genetic correlation between the particular trait and all-cause mortality.⁶⁹ Further research on the genetic determinants of disease prognosis and survival may also increase the utility of PRS in understanding mortality risk.

There is the potential for misuse of polygenic risk scores, including the composite PRS we have developed for predicting mortality. In particular, we urge great caution in the deployment of PRS and advocate for the creation of appropriate measures in order to prevent the misuse of PRS, e.g., for embryo screening⁷⁰ or in ways that could put individuals at risk of “genetic discrimination.”⁷¹ As PRS for predicting various outcomes become increasingly available, a suitable regulatory framework for implementation will be needed to allow for the utilization of PRS to improve healthcare while protecting individuals from harm due to potential misuse.

In conclusion, our results suggest that by combining knowledge gained from GWASs of complex traits, it may be possible to identify individuals who are expected to live substantially longer or shorter. In light of the ethical repercussions of using genetics to make predictions regarding an individual’s life course prior to or at birth, we argue that the cPRS may be most useful for counseling those in early adulthood about their genetic risk. In partic-

ular, the results of our analysis highlight the importance of considering genetic risk in the context of clinical risk factors measured in adulthood; thus, the cPRS may be useful in advising patients on the importance of certain lifestyle choices associated with mortality risk. Using the cPRS in this way would require validation of the cPRS outside of the UK Biobank.

Data and Code Availability

Data from the UK Biobank are available by application to the UK Biobank. The data needed to generate the trait-specific PRS (i.e., RS IDs and SNP weights) used in the main analysis and code to construct and evaluate the composite PRS are available on GitHub (see [CompositePRS Code in Web Resources](#)). The GWAS summary statistic data used to construct the LDpred PRS for breast cancer and stroke can be downloaded from the BCAC and Megastroke sites, respectively (see [Web Resources](#)).

Supplemental Data

Supplemental Data can be found online at <https://doi.org/10.1016/j.ajhg.2020.07.002>.

Acknowledgments

This research has been conducted using the UK Biobank Resource under Application Number 17712. This research was supported by the Patient-Centered Outcomes Research Institute (ME-1602-34530), the National Institutes of Health (1 R01 HG010480-01), and the Intramural Research Program, Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health.

Declaration of Interests

The authors declare no competing interests.

Received: April 14, 2020

Accepted: July 1, 2020

Published: August 5, 2020

Web Resources

BCAC GWAS Summary Results, <http://bcac.ccge.medschl.cam.ac.uk/bcacdata/oncoarray/oncoarray-and-combined-summary-result/gwas-summary-results-breast-cancer-risk-2017/>

CDC WONDER Database: Underlying Cause of Death, <https://wonder.cdc.gov/ucd-icd10.html>

CompositePRS Code, <https://github.com/meisnera/CompositePRS>

Megastroke, <http://www.megastroke.org/index.html>

Neale Lab, <http://www.nealelab.is/uk-biobank>

Office of National Statistics, <https://www.nomisweb.co.uk/datasets/mortsa>

R Software, <https://www.r-project.org/>

rms Package, <https://cran.r-project.org/package=rms>

survival Package, <https://cran.r-project.org/package=survival>

References

1. Visscher, P.M., Wray, N.R., Zhang, Q., Sklar, P., McCarthy, M.I., Brown, M.A., and Yang, J. (2017). 10 years of GWAS discovery: biology, function, and translation. *Am. J. Hum. Genet.* *101*, 5–22.
2. Zeng, J., de Vlaming, R., Wu, Y., Robinson, M.R., Lloyd-Jones, L.R., Yengo, L., Yap, C.X., Xue, A., Sidorenko, J., McRae, A.F., et al. (2018). Signatures of negative selection in the genetic architecture of human complex traits. *Nat. Genet.* *50*, 746–753.
3. Zhang, Y., Qi, G., Park, J.H., and Chatterjee, N. (2018). Estimation of complex effect-size distributions using summary-level statistics from genome-wide association studies across 32 complex traits. *Nat. Genet.* *50*, 1318–1326.
4. Khera, A.V., Emdin, C.A., Drake, I., Natarajan, P., Bick, A.G., Cook, N.R., Chasman, D.I., Baber, U., Mehran, R., Rader, D.J., et al. (2016). Genetic risk, adherence to a healthy lifestyle, and coronary disease. *N. Engl. J. Med.* *375*, 2349–2358.
5. Lewis, C.M., and Vassos, E. (2017). Prospects for using risk scores in polygenic medicine. *Genome Med.* *9*, 96.
6. Garcia-Closas, M., Rothman, N., Figueroa, J.D., Prokunina-Olsson, L., Han, S.S., Baris, D., Jacobs, E.J., Malats, N., De Vivo, I., Albanes, D., et al. (2013). Common genetic polymorphisms modify the effect of smoking on absolute risk of bladder cancer. *Cancer Res.* *73*, 2211–2220.
7. Chatterjee, N., Shi, J., and García-Closas, M. (2016). Developing and evaluating polygenic risk prediction models for stratified disease prevention. *Nat. Rev. Genet.* *17*, 392–406.
8. Maas, P., Barrdahl, M., Joshi, A.D., Auer, P.L., Gaudet, M.M., Milne, R.L., Schumacher, F.R., Anderson, W.F., Check, D., Chattopadhyay, S., et al. (2016). Breast cancer risk from modifiable and nonmodifiable risk factors among white women in the United States. *JAMA Oncol.* *2*, 1295–1302.
9. Frampton, M.J.E., Law, P., Litchfield, K., Morris, E.J., Kerr, D., Turnbull, C., Tomlinson, I.P., and Houlston, R.S. (2016). Implications of polygenic risk for personalised colorectal cancer screening. *Ann. Oncol.* *27*, 429–434.
10. Seibert, T.M., Fan, C.C., Wang, Y., Zuber, V., Karunamuni, R., Parsons, J.K., Eeles, R.A., Easton, D.F., Kote-Jarai, Z., Al Olama, A.A., et al.; PRACTICAL Consortium* (2018). Polygenic hazard score to guide screening for aggressive prostate cancer: development and validation in large scale cohorts. *BMJ* *360*, j5757.
11. Mavaddat, N., Pharoah, P.D.P., Michailidou, K., Tyrer, J., Brook, M.N., Bolla, M.K., Wang, Q., Dennis, J., Dunning, A.M., Shah, M., et al. (2015). Prediction of breast cancer risk based on profiling with common genetic variants. *J. Natl. Cancer Inst.* *107*, djv036.
12. Hsu, L., Jeon, J., Brenner, H., Gruber, S.B., Schoen, R.E., Berndt, S.I., Chan, A.T., Chang-Claude, J., Du, M., Gong, J., et al.; Colorectal Transdisciplinary (CORECT) Study; and Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO) (2015). A model to determine colorectal cancer risk using common genetic susceptibility loci. *Gastroenterology* *148*, 1330–9.e14.
13. Mega, J.L., Stitzel, N.O., Smith, J.G., Chasman, D.I., Caulfield, M., Devlin, J.J., Nordio, F., Hyde, C., Cannon, C.P., Sacks, F., et al. (2015). Genetic risk, coronary heart disease events, and the clinical benefit of statin therapy: an analysis of primary and secondary prevention trials. *Lancet* *385*, 2264–2271.
14. Natarajan, P., Young, R., Stitzel, N.O., Padmanabhan, S., Baber, U., Mehran, R., Sartori, S., Fuster, V., Reilly, D.F., Butterworth, A., et al. (2017). Polygenic risk score identifies subgroup with higher burden of atherosclerosis and greater relative benefit from statin therapy in the primary prevention setting. *Circulation* *135*, 2091–2101.
15. Mavaddat, N., Michailidou, K., Dennis, J., Lush, M., Fachal, L., Lee, A., Tyrer, J.P., Chen, T.H., Wang, Q., Bolla, M.K., et al.; ABCTB Investigators; kConFab/AOCS Investigators; and NBCS Collaborators (2019). Polygenic risk scores for prediction of breast cancer and breast cancer subtypes. *Am. J. Hum. Genet.* *104*, 21–34.
16. Torkamani, A., Wineinger, N.E., and Topol, E.J. (2018). The personal and clinical utility of polygenic risk scores. *Nat. Rev. Genet.* *19*, 581–590.
17. Lambert, S.A., Abraham, G., and Inouye, M. (2019). Towards clinical utility of polygenic risk scores. *Hum. Mol. Genet.* *28* (R2), R133–R142.
18. Wald, N.J., and Old, R. (2019). The illusion of polygenic disease risk prediction. *Genet. Med.* *21*, 1705–1707.
19. Khera, A.V., Chaffin, M., Wade, K.H., Zahid, S., Brancale, J., Xia, R., Distefano, M., Senol-Cosar, O., Haas, M.E., Bick, A., et al. (2019). Polygenic prediction of weight and obesity trajectories from birth to adulthood. *Cell* *177*, 587–596.e9.
20. Langenberg, C., Sharp, S.J., Franks, P.W., Scott, R.A., Deloukas, P., Forouhi, N.G., Froguel, P., Groop, L.C., Hansen, T., Palla, L., et al. (2014). Gene-lifestyle interaction and type 2 diabetes: the EPIC interact case-cohort study. *PLoS Med.* *11*, e1001647.
21. Rudolph, A., Song, M., Brook, M.N., Milne, R.L., Mavaddat, N., Michailidou, K., Bolla, M.K., Wang, Q., Dennis, J., Wilcox, A.N., et al. (2018). Joint associations of a polygenic risk score and environmental risk factors for breast cancer in the Breast Cancer Association Consortium. *Int. J. Epidemiol.* *47*, 526–536.
22. Hedley, A.J., Wong, C.M., Thach, T.Q., Ma, S., Lam, T.H., and Anderson, H.R. (2002). Cardiorespiratory and all-cause mortality after restrictions on sulphur content of fuel in Hong Kong: an intervention study. *Lancet* *360*, 1646–1652.
23. Anthonisen, N.R., Skeans, M.A., Wise, R.A., Manfreda, J., Kanner, R.E., Connett, J.E.; and Lung Health Study Research Group (2005). The effects of a smoking cessation intervention on 14.5-year mortality: a randomized clinical trial. *Ann. Intern. Med.* *142*, 233–239.
24. Grooteman, M.P.C., van den Dorpel, M.A., Bots, M.L., Penne, E.L., van der Weerd, N.C., Mazairac, A.H.A., den Hoedt, C.H., van der Tweel, I., Lévesque, R., Nubé, M.J., et al.; CONTRAST Investigators (2012). Effect of online hemodiafiltration on all-cause mortality and cardiovascular outcomes. *J. Am. Soc. Nephrol.* *23*, 1087–1096.
25. Mohiuddin, S.M., Mooss, A.N., Hunter, C.B., Grollmes, T.L., Cloutier, D.A., and Hilleman, D.E. (2007). Intensive smoking cessation intervention reduces mortality in high-risk smokers with cardiovascular disease. *Chest* *131*, 446–452.
26. Timmers, P.R., Mounier, N., Lall, K., Fischer, K., Ning, Z., Feng, X., Bretherick, A.D., Clark, D.W., Agbessi, M., Ahsan, H., et al.; eQTLGen Consortium (2019). Genomics of 1 million parent lifespans implicates novel pathways and common diseases and distinguishes survival chances. *eLife* *8*, 8.
27. Wright, K.M., Rand, K.A., Kermany, A., Noto, K., Curtis, D., Garrigan, D., Slinkov, D., Dorfman, I., Granka, J.M., Byrnes, J., et al. (2019). A prospective analysis of genetic variants associated with human lifespan. *G3 (Bethesda)* *9*, 2863–2878.
28. Melzer, D., Pilling, L.C., and Ferrucci, L. (2020). The genetics of human ageing. *Nat. Rev. Genet.* *21*, 88–101.

29. Public Health England (2017). Health Profile for England:2017. Chapter 2: Major causes of death and how they have changed. <https://www.gov.uk/government/publications/health-profile-for-england/chapter-2-major-causes-of-death-and-how-they-have-changed>.
30. Ganna, A., Rivadeneira, F., Hofman, A., Uitterlinden, A.G., Magnusson, P.K.E., Pedersen, N.L., Ingelsson, E., and Tiemeier, H. (2013). Genetic determinants of mortality. Can findings from genome-wide association studies explain variation in human mortality? *Hum. Genet.* *132*, 553–561.
31. Joshi, P.K., Pirastu, N., Kentistou, K.A., Fischer, K., Hofer, E., Schraut, K.E., Clark, D.W., Nutile, T., Barnes, C.L.K., Timmers, P.R.H.J., et al. (2017). Genome-wide meta-analysis associates HLA-DQA1/DRB1 and LPA and lifestyle factors with human longevity. *Nat. Commun.* *8*, 910.
32. da Silva, A.A., de Mello, R.G.B., Schaan, C.W., Fuchs, F.D., Redline, S., and Fuchs, S.C. (2016). Sleep duration and mortality in the elderly: a systematic review with meta-analysis. *BMJ Open* *6*, e008119.
33. Cappuccio, F.P., D'Elia, L., Strazzullo, P., and Miller, M.A. (2010). Sleep duration and all-cause mortality: a systematic review and meta-analysis of prospective studies. *Sleep* *33*, 585–592.
34. Liu, T.Z., Xu, C., Rota, M., Cai, H., Zhang, C., Shi, M.J., Yuan, R.X., Weng, H., Meng, X.Y., Kwong, J.S.W., and Sun, X. (2017). Sleep duration and risk of all-cause mortality: A flexible, non-linear, meta-regression of 40 prospective cohort studies. *Sleep Med. Rev.* *32*, 28–36.
35. Buniello, A., MacArthur, J.A.L., Cerezo, M., Harris, L.W., Hayhurst, J., Malangone, C., McMahon, A., Morales, J., Mountjoy, E., Sollis, E., et al. (2019). The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* *47* (D1), D1005–D1012.
36. Hunt, S.E., McLaren, W., Gil, L., Thormann, A., Schuilenburg, H., Sheppard, D., Parton, A., Armean, I.M., Trevanion, S.J., Flicek, P., et al. (2018). Ensembl variation resources. *Database*. 2018. <https://doi.org/10.1093/database/bay119>.
37. Rutten-Jacobs, L.C.A., Larsson, S.C., Malik, R., Rannikmäe, K., Sudlow, C.L., Dichgans, M., Markus, H.S., Traylor, M.; MEGASTROKE consortium; and International Stroke Genetics Consortium (2018). Genetic risk, incident stroke, and the benefits of adhering to a healthy lifestyle: cohort study of 306c473 UK Biobank participants. *BMJ* *363*, k4168.
38. Liu, M., Jiang, Y., Wedow, R., Li, Y., Brazel, D.M., Chen, F., Datta, G., Davila-Velderrain, J., McGuire, D., Tian, C., et al.; 23andMe Research Team; and HUNT All-In Psychiatry (2019). Association studies of up to 1.2 million individuals yield new insights into the genetic etiology of tobacco and alcohol use. *Nat. Genet.* *51*, 237–244.
39. Yengo, L., Sidorenko, J., Kemper, K.E., Zheng, Z., Wood, A.R., Weedon, M.N., Frayling, T.M., Hirschhorn, J., Yang, J., Visscher, P.M.; and GIANT Consortium (2018). Meta-analysis of genome-wide association studies for height and body mass index in ~700000 individuals of European ancestry. *Hum. Mol. Genet.* *27*, 3641–3649.
40. Willer, C.J., Schmidt, E.M., Sengupta, S., Peloso, G.M., Gustafsson, S., Kanoni, S., Ganna, A., Chen, J., Buchkovich, M.L., Mora, S., et al.; Global Lipids Genetics Consortium (2013). Discovery and refinement of loci associated with lipid levels. *Nat. Genet.* *45*, 1274–1283.
41. Scott, R.A., Lagou, V., Welch, R.P., Wheeler, E., Montasser, M.E., Luan, J., Mägi, R., Strawbridge, R.J., Rehnberg, E., Gustafsson, S., et al.; DIABetes Genetics Replication and Meta-analysis (DIAGRAM) Consortium (2012). Large-scale association analyses identify new loci influencing glycemic traits and provide insight into the underlying biological pathways. *Nat. Genet.* *44*, 991–1005.
42. Li, M., Li, Y., Weeks, O., Mijatovic, V., Teumer, A., Huffman, J.E., Tromp, G., Fuchsberger, C., Gorski, M., Lytykäinen, L.P., et al.; CHARGE Glycemic-T2D Working Group; and CHARGE Blood Pressure Working Group (2017). SOS2 and ACP1 loci identified through large-scale exome chip analysis regulate kidney development and function. *J. Am. Soc. Nephrol.* *28*, 981–994.
43. Dashti, H.S., Jones, S.E., Wood, A.R., Lane, J.M., van Hees, V.T., Wang, H., Rhodes, J.A., Song, Y., Patel, K., Anderson, S.G., et al. (2019). Genome-wide association study identifies genetic loci for self-reported habitual sleep duration supported by accelerometer-derived estimates. *Nat. Commun.* *10*, 1100.
44. Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., et al. (2015). UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* *12*, e1001779.
45. Eastwood, S.V., Mathur, R., Atkinson, M., Brophy, S., Sudlow, C., Flaig, R., de Lusignan, S., Allen, N., and Chaturvedi, N. (2016). Algorithms for the capture and adjudication of prevalent and incident diabetes in UK Biobank. *PLoS ONE* *11*, e0162388.
46. Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A., Bender, D., Maller, J., Sklar, P., de Bakker, P.I., Daly, M.J., and Sham, P.C. (2007). PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* *81*, 559–575.
47. Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., Korbel, J.O., Marchini, J.L., McCarthy, S., McVean, G.A., Abecasis, G.R.; and 1000 Genomes Project Consortium (2015). A global reference for human genetic variation. *Nature* *526*, 68–74.
48. Zou, G. (2004). A modified poisson regression approach to prospective studies with binary data. *Am. J. Epidemiol.* *159*, 702–706.
49. Levey, A.S., de Jong, P.E., Coresh, J., El Nahas, M., Astor, B.C., Matsushita, K., Gansevoort, R.T., Kasiske, B.L., and Eckardt, K.U. (2011). The definition, classification, and prognosis of chronic kidney disease: a KDIGO Controversies Conference report. *Kidney Int.* *80*, 17–28.
50. Harrell, F.E., Jr., Califf, R.M., Pryor, D.B., Lee, K.L., and Rosati, R.A. (1982). Evaluating the yield of medical tests. *JAMA* *247*, 2543–2546.
51. Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis* (New York: Springer-Verlag).
52. Zeileis, A. (2004). Econometric computing with HC and HAC covariance matrix estimators. *J. Stat. Softw.* *11*, 1–17.
53. Zeileis, A. (2006). Object-oriented computation of sandwich estimators. *J. Stat. Softw.* *16*, 1–16.
54. Mostafavi, H., Berisa, T., Day, F.R., Perry, J.R.B., Przeworski, M., and Pickrell, J.K. (2017). Identifying genetic variants that affect viability in large cohorts. *PLoS Biol.* *15*, e2002458.
55. Pilling, L.C., Atkins, J.L., Bowman, K., Jones, S.E., Tyrrell, J., Beaumont, R.N., Ruth, K.S., Tuke, M.A., Yaghootkar, H., Wood, A.R., et al. (2016). Human longevity is influenced by

- many genetic variants: evidence from 75,000 UK Biobank participants. *Aging (Albany N.Y.)* 8, 547–560.
56. Pilling, L.C., Kuo, C.-L., Sicinski, K., Tamosauskaite, J., Kuchel, G.A., Harries, L.W., Herd, P., Wallace, R., Ferrucci, L., and Melzer, D. (2017). Human longevity: 25 genetic loci associated in 389,166 UK biobank participants. *Aging (Albany N.Y.)* 9, 2504–2520.
 57. Marioni, R.E., Ritchie, S.J., Joshi, P.K., Hagenaars, S.P., Okbay, A., Fischer, K., Adams, M.J., Hill, W.D., Davies, G., Nagy, R., et al.; Social Science Genetic Association Consortium (2016). Genetic variants linked to education predict longevity. *Proc. Natl. Acad. Sci. USA* 113, 13366–13371.
 58. Dziedzic, T., Slowik, A., Gryz, E.A., and Szczudlik, A. (2004). Lower serum triglyceride level is associated with increased stroke severity. *Stroke* 35, e151–e152.
 59. Jain, M., Jain, A., Yerragondlu, N., Brown, R.D., Rabinstein, A., Jahromi, B.S., Vaidyanathan, L., Blyth, B., and Stead, L.G. (2013). The triglyceride paradox in stroke survivors: A prospective study. *Neurosci. J.* 2013, 870608.
 60. Ryu, W.S., Lee, S.H., Kim, C.K., Kim, B.J., and Yoon, B.W. (2010). Effects of low serum triglyceride on stroke mortality: a prospective follow-up study. *Atherosclerosis* 212, 299–304.
 61. Li, W., Liu, M., Wu, B., Liu, H., Wang, L.C., and Tan, S. (2008). Serum lipid levels and 3-month prognosis in Chinese patients with acute stroke. *Adv. Ther.* 25, 329–341.
 62. Beekman, M., Blanché, H., Perola, M., Hervonen, A., Bezrukov, V., Sikora, E., Flachsbarth, F., Christiansen, L., De Craen, A.J.M., Kirkwood, T.B.L., et al.; GEHA consortium (2013). Genome-wide linkage analysis for human longevity: Genetics of Healthy Aging Study. *Aging Cell* 12, 184–193.
 63. Joshi, P.K., Fischer, K., Schraut, K.E., Campbell, H., Esko, T., and Wilson, J.F. (2016). Variants near CHRNA3/5 and APOE have age- and sex-related effects on human lifespan. *Nat. Commun.* 7, 11174.
 64. Sakaue, S., Kanai, M., Karjalainen, J., Akiyama, M., Kurki, M., Matoba, N., Takahashi, A., Hirata, M., Kubo, M., Matsuda, K., et al.; FinnGen (2020). Trans-biobank analysis with 676,000 individuals elucidates the association of polygenic risk scores of complex traits with human lifespan. *Nat. Med.* 26, 542–548.
 65. Fry, A., Littlejohns, T.J., Sudlow, C., Doherty, N., Adamska, L., Sprosen, T., Collins, R., and Allen, N.E. (2017). Comparison of sociodemographic and health-related characteristics of UK Biobank participants with those of the general population. *Am. J. Epidemiol.* 186, 1026–1034.
 66. Keyes, K.M., and Westreich, D. (2019). UK Biobank, big data, and the consequences of non-representativeness. *Lancet* 393, 1297.
 67. Privé, F., Vilhjálmsón, B.J., Aschard, H., and Blum, M.G.B. (2019). Making the most of clumping and thresholding for polygenic scores. *Am. J. Hum. Genet.* 105, 1213–1221.
 68. Vilhjálmsón, B.J., Yang, J., Finucane, H.K., Gusev, A., Lindström, S., Ripke, S., Genovese, G., Loh, P.R., Bhatia, G., Do, R., et al.; Schizophrenia Working Group of the Psychiatric Genomics Consortium, Discovery, Biology, and Risk of Inherited Variants in Breast Cancer (DRIVE) study (2015). Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *Am. J. Hum. Genet.* 97, 576–592.
 69. Krapohl, E., Patel, H., Newhouse, S., Curtis, C.J., von Stumm, S., Dale, P.S., Zabaneh, D., Breen, G., O'Reilly, P.F., and Plomin, R. (2018). Multi-polygenic score approach to trait prediction. *Mol. Psychiatry* 23, 1368–1374.
 70. Treff, N.R., Eccles, J., Lello, L., Bechor, E., Hsu, J., Plunkett, K., Zimmerman, R., Rana, B., Samoilenko, A., Hsu, S., and Tellier, L.C.A.M. (2019). Utility and first clinical application of screening embryos for polygenic disease risk reduction. *Front. Endocrinol. (Lausanne)* 10, 845.
 71. Godard, B., Raeburn, S., Pembrey, M., Bobrow, M., Farndon, P., and Aymé, S. (2003). Genetic information and testing in insurance and employment: technical, social and ethical issues. *Eur. J. Hum. Genet.* 11 (Suppl 2), S123–S142.

The American Journal of Human Genetics, Volume 107

Supplemental Data

Combined Utility of 25 Disease and Risk Factor

Polygenic Risk Scores

for Stratifying Risk of All-Cause Mortality

Allison Meisner, Prosenjit Kundu, Yan Dora Zhang, Lauren V. Lan, Sungwon Kim, Disha Ghandwani, Parichoy Pal Choudhury, Sonja I. Berndt, Neal D. Freedman, Montserrat Garcia-Closas, and Nilanjan Chatterjee

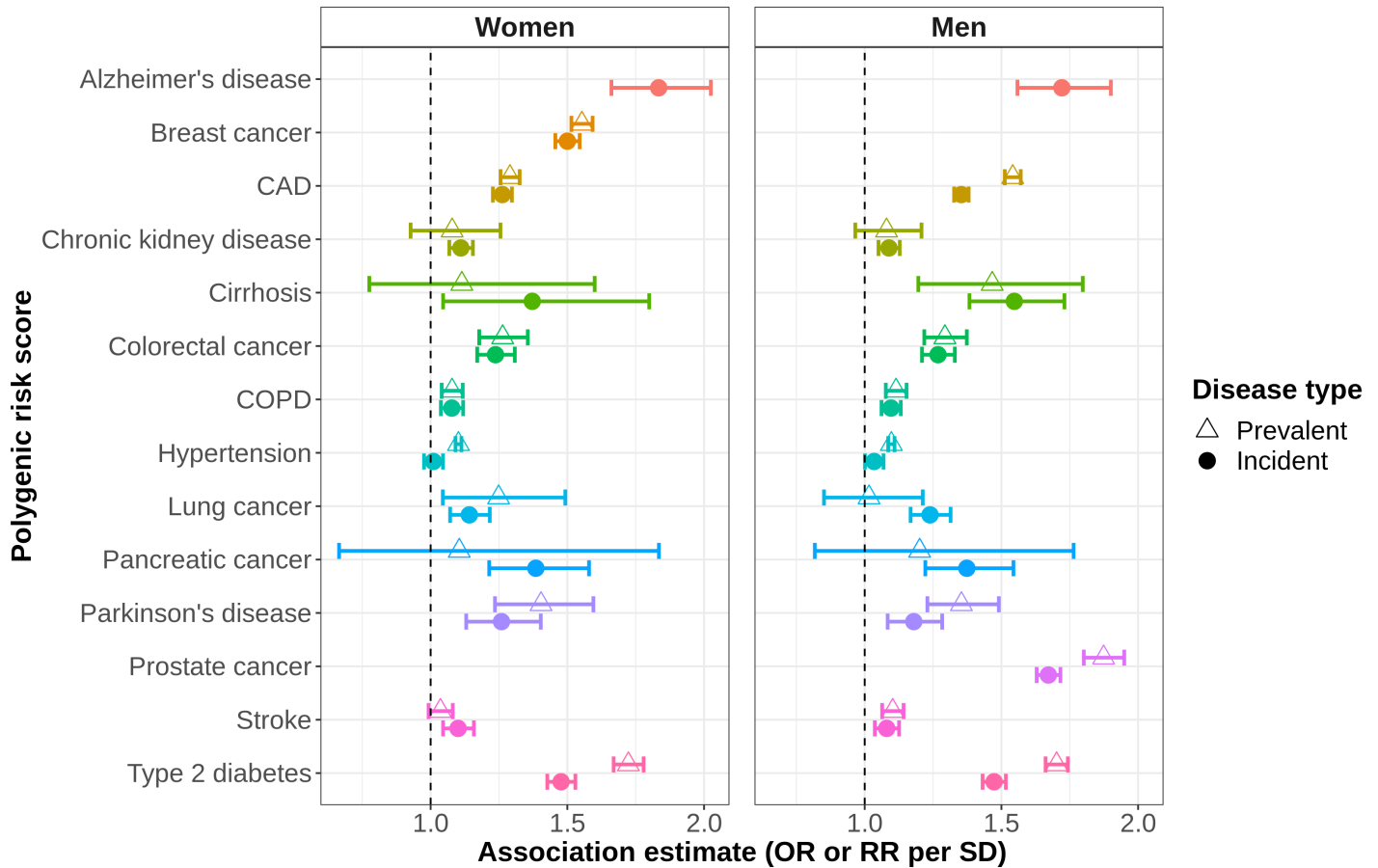


Figure S1. The estimated association between each disease PRS and prevalent and incident disease. The results are presented for women (left panel) and men (right panel) separately. For prevalent disease (open triangles in the plot), sex-specific logistic regression models were fit in the full cohort. For incident disease (closed circles in the plot), sex-specific modified Poisson regression models with robust standard error estimates were fit to the full cohort, excluding individuals with the disease at baseline (prevalent cases). All models included adjustment for age at entry. The estimates are presented as the estimated OR or RR per standard deviation of the PRS. The horizontal lines indicate 95% confidence intervals. As the number of prevalent cases of Alzheimer's disease was quite low for both men and women, these estimates are not presented. CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; OR: odds ratio; RR: relative risk; SD: standard deviation; PRS: polygenic risk score.

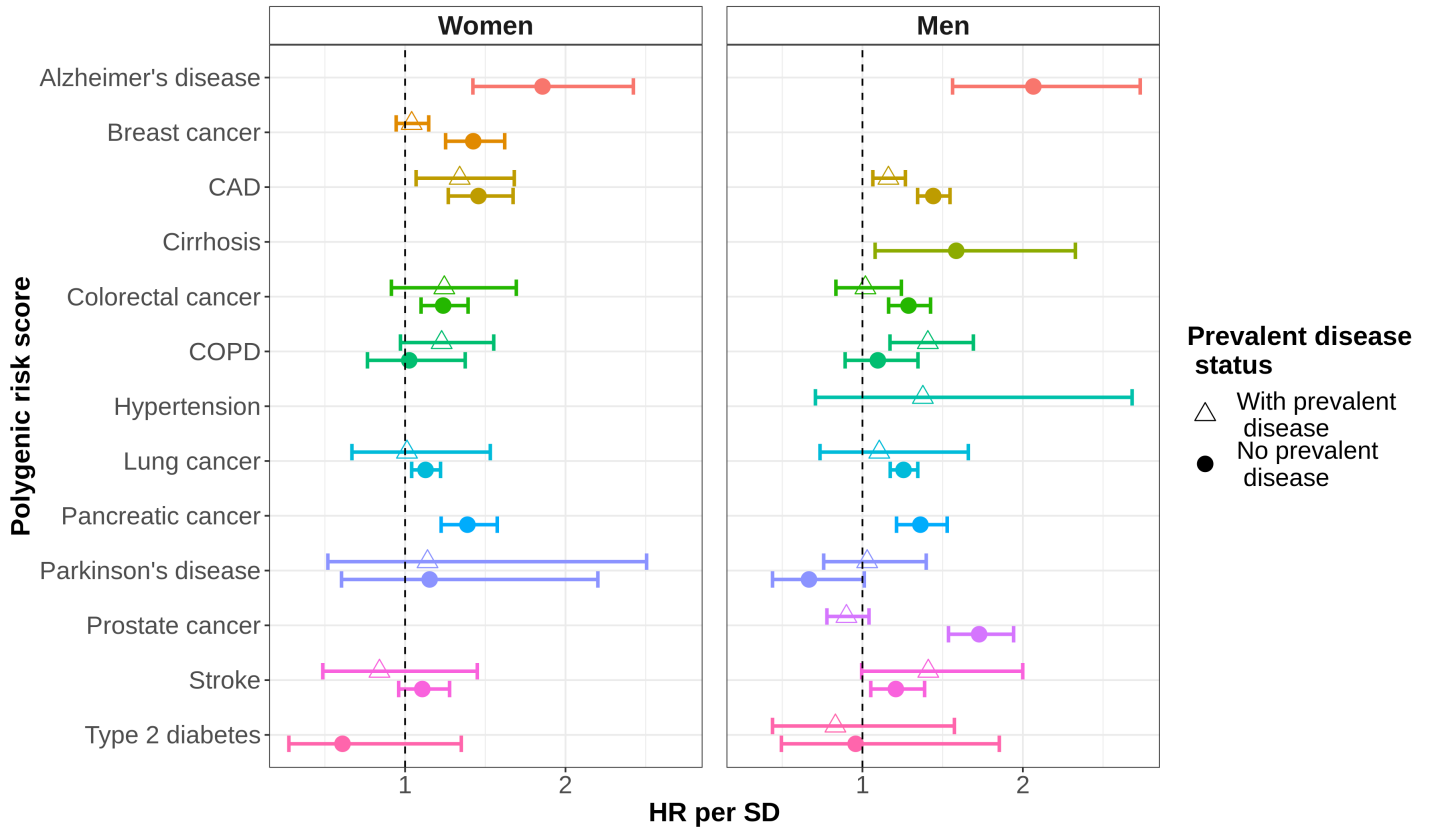


Figure S2. Cause-specific mortality results, stratified by the presence of disease at study baseline. For each disease, we used the data from the full cohort to evaluate the association between the disease PRS and mortality from the disease based on sex-specific Cox proportional hazards models of age at death in individuals with the disease at baseline (open triangles in the plot) and in individuals without the disease at baseline (closed circles in the plot). Deaths from other causes were treated as censoring events. Some causes did not have enough observations or deaths to yield stable estimates (< 30 observations or < 6 deaths); in these cases, estimates are not provided. Each PRS was standardized to have unit variance so the estimates correspond to the HR per SD of the PRS. The horizontal lines indicate 95% confidence intervals. CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; HR: hazard ratio; SD: standard deviation; PRS: polygenic risk score.

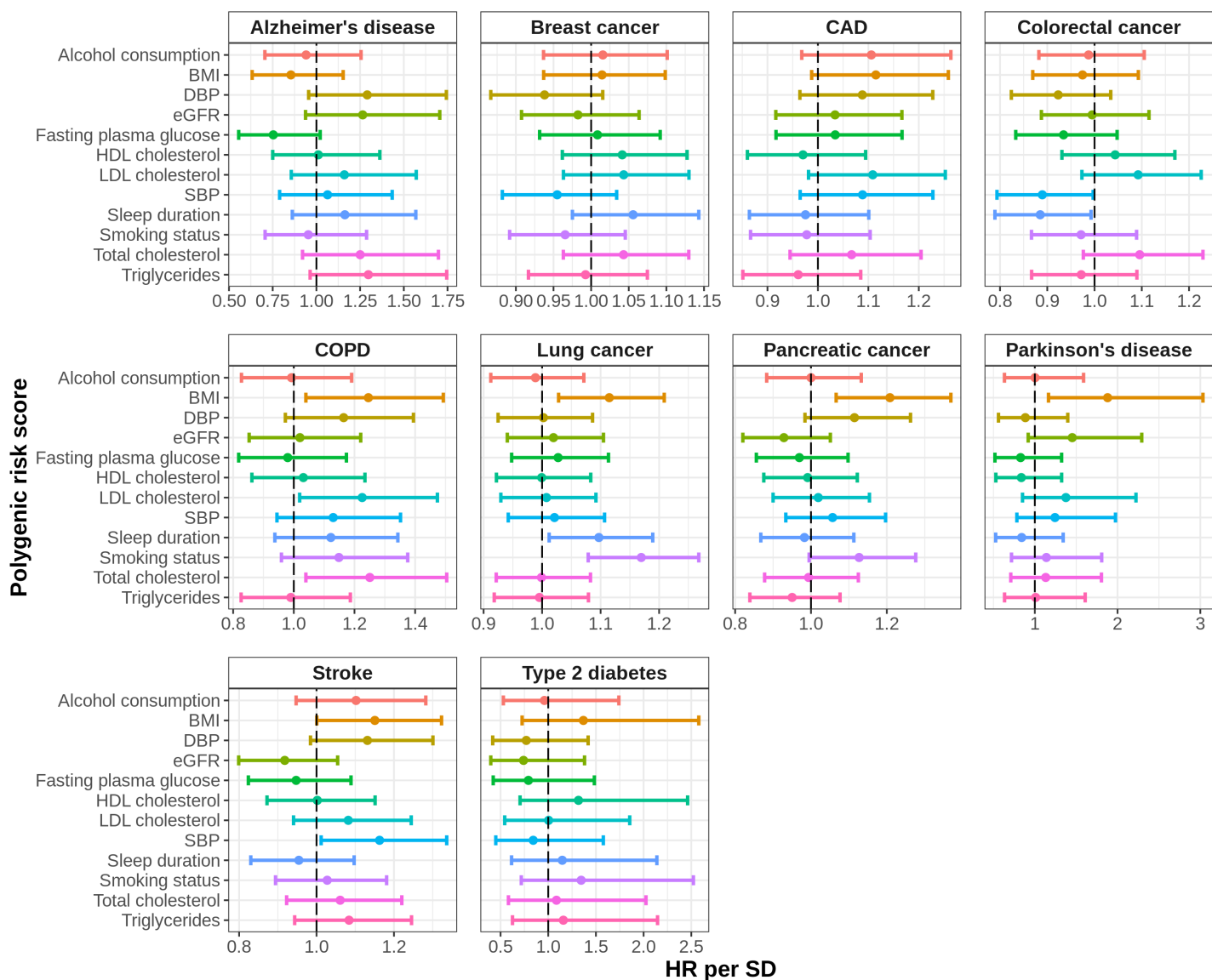


Figure S3. The estimated association between each mortality risk factor PRS and mortality due to each of the top causes of death among women. For each disease, we evaluated the association between each of the risk factor PRS and mortality from the disease based on Cox proportional hazards models of age at death in women in the full cohort. Deaths from other causes were treated as censoring events. Some causes did not have enough deaths to yield stable estimates (< 6 deaths); in these cases, estimates are not provided. Each PRS was standardized to have unit variance so the estimates correspond to the HR per SD of the PRS. The horizontal lines indicate 95% confidence intervals. BMI: body mass index; CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; HDL: high-density lipoprotein; LDL: low-density lipoprotein; SBP: systolic blood pressure; HR: hazard ratio; SD: standard deviation; PRS: polygenic risk score.

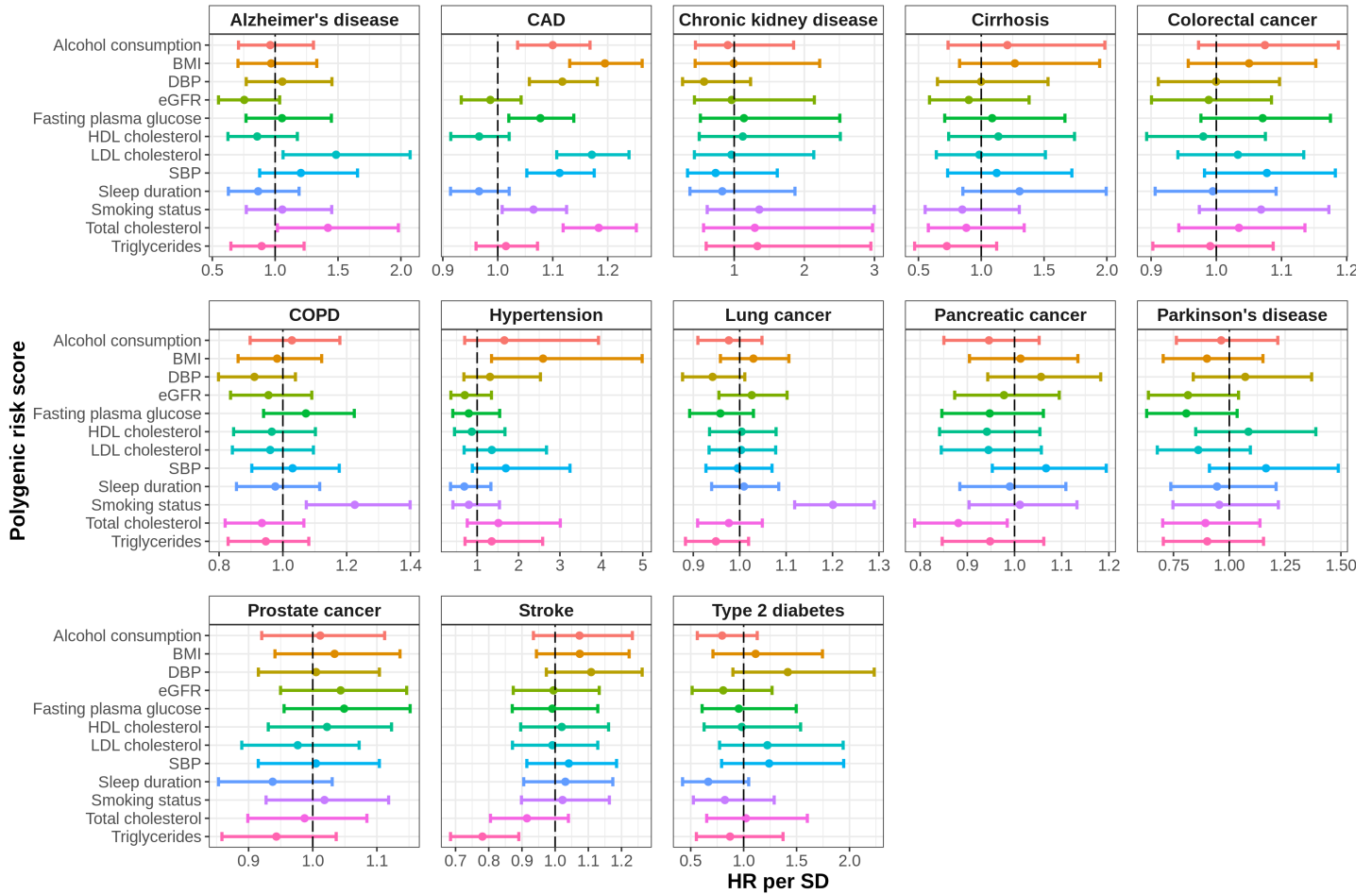


Figure S4. The estimated association between each mortality risk factor PRS and mortality due to each of the top causes of death among men. For each disease, we evaluated the association between each of the risk factor PRS and mortality from the disease based on Cox proportional hazards models of age at death in men in the full cohort. Deaths from other causes were treated as censoring events. Each PRS was standardized to have unit variance so the estimates correspond to the HR per SD of the PRS. The horizontal lines indicate 95% confidence intervals. BMI: body mass index; CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; HDL: high-density lipoprotein; LDL: low-density lipoprotein; SBP: systolic blood pressure; HR: hazard ratio; SD: standard deviation; PRS: polygenic risk score.

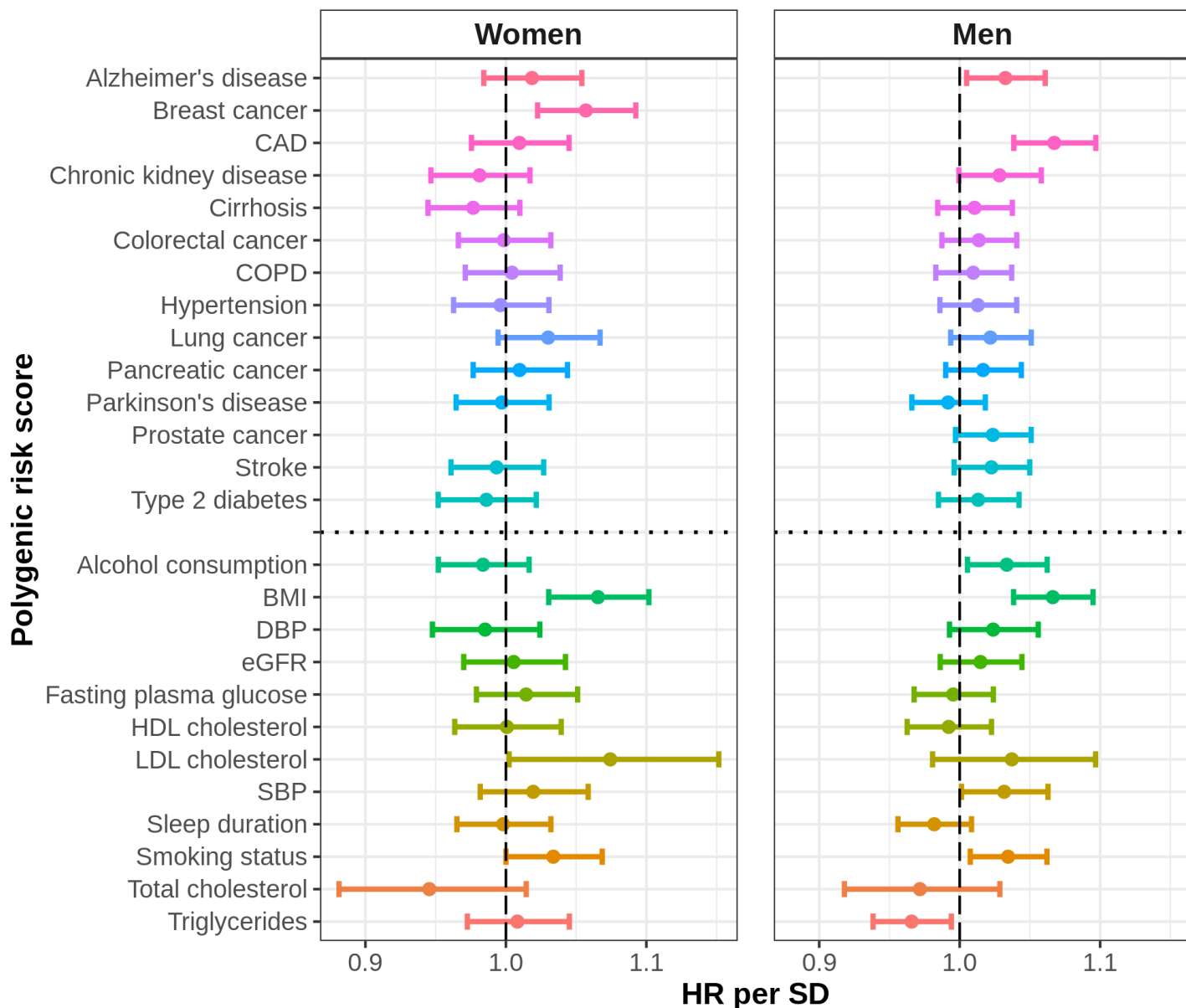


Figure S5. Association of trait-specific PRS with all-cause mortality in the training data based on models with all 25 PRS. The estimates are based on sex-specific Cox proportional hazards models of age at death with all 25 PRS, fit in the training data. These association estimates were used to weight each PRS to form the cPRS. Each PRS was standardized to have unit variance so the estimates correspond to the HR per SD of the PRS. The horizontal lines indicate 95% confidence intervals. BMI: body mass index; CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; HDL: high-density lipoprotein; LDL: low-density lipoprotein; SBP: systolic blood pressure; HR: hazard ratio; SD: standard deviation; PRS: polygenic risk score.

Table S1. ICD-10 codes for the top causes of death. The top causes of death (“CDC Definition”) based on the CDC WONDER database and the corresponding specific cause of death included in the analysis (“Our Definition”) are both presented. Ranking in the US based on data for 2017 from CDC WONDER for non-Hispanic whites aged 40 and over; ranking in the UK based on data for 2017 from the Office of National Statistics for individuals aged 40 and over.

Ranking in US (UK)		CDC Definition		Our Definition	
Women	Men	Cause	ICD-10 codes	Cause	ICD-10 codes
1 (2)	1 (2)	Diseases of heart	I00-I09, I11, I13, I20-I51	CAD	I20-I25
2 (1)	2 (1)	Malignant neoplasms	C00-C97	Pancreatic	C25
				Colorectal	C18-C20
				Breast	C50
				Lung	C33-C34
		Prostate	C61		
3 (4)	3 (3)	Chronic lower respiratory diseases	J40-J47	Chronic obstructive pulmonary disease	J41-J44
4 (5)	5 (5)	Alzheimer’s disease	G30	Alzheimer’s disease	G30
5 (3)	4 (4)	Cerebrovascular diseases	I60-I69	Stroke	I60, I61, I63, I64
6 (6)	6 (9)	Diabetes mellitus	E10-E14	Type 2 diabetes	E11
7 (10)	8 (10)	Nephritis, nephrotic syndrome and nephrosis	N00-N07, N17-N19, N25-N27	Chronic kidney disease	N18
8 (11)	10 (11)	Essential hypertension and hypertensive renal disease	I10, I12, I15	Hypertension	I10
9 (7)	7 (6)	Chronic liver disease and cirrhosis	K70, K73-K74	Alcoholic liver cirrhosis	K70.3
10 (8)	9 (7)	Parkinson’s disease	G20-G21	Parkinson’s disease	G20

CAD: Coronary artery disease; CDC: Centers for Disease Control; ICD: International Classification of Diseases; WONDER: Wide-ranging ONline Data for Epidemiologic Research.

Table S2. Methods for identifying prevalent and incident cases of each disease included in the analysis.

Cause of death	ICD Codes		Prevalent Definition	Incident Definition
	ICD9	ICD10		
Coronary artery disease	410-414	I20-I25	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment (b) Self-report: self-reported CAD at baseline	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death
Pancreatic cancer	157	C25	Cancer registry: one of the ICD (9/10) codes in the cancer registry with an initial date prior to the date of baseline assessment	Cancer registry: one of the ICD (9/10) codes in the cancer registry with an initial date after date of baseline assessment
Colorectal cancer	153, 154.0, 154.1, 154.8	C18-C20		
Breast cancer	174	C50		
Lung cancer	162	C33-C34		
Prostate cancer	185	C61		
COPD	491, 492, 496	J41-J44	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment (b) Self-report: self-reported COPD, emphysema, or chronic bronchitis at baseline	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death
Alzheimer's disease	331.0	G30 and F00	(a) HES: one of the ICD (9/10) codes in the primary or any secondary position with an initial code date is prior to the date of baseline assessment.	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death
Stroke	430, 431, 434, 436	I60, I61, I63, I64	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment (b) Self-report: self-reported stroke at baseline	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death
Type 2 diabetes	Defined based on algorithms in Eastwood et al. (1)			
Chronic kidney disease	585	N18	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death
Hypertension	401	I10	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment (b) Self-report: (i) self-reported essential hypertension or "any hypertension" but not "gestational hypertension/pre-eclampsia" at	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death

			baseline or (ii) hypertension medication usage at baseline (c) SBP/DBP measures: systolic blood pressure ≥ 140 mmHg, or diastolic blood pressure ≥ 90 mmHg at baseline	
Alcoholic liver cirrhosis	571.2	K70.3	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death
Parkinson's disease	332.0	G20	(a) HES: ICD (9/10) codes in HES in the primary or secondary position with an initial code date prior to the date of baseline assessment (b) Self-report: self-reported Parkinson's disease at baseline	(a) HES: one of the ICD (9/10) codes in HES in the primary or secondary position with an initial code date after the date of baseline assessment (b) Mortality data: one of the ICD10 codes listed as a primary or secondary cause of death

COPD: chronic obstructive pulmonary disease; HES: hospital episode statistics data; ICD: International Classification of Diseases.

Table S3. Conversion of self-reported alcohol intake to grams of alcohol per day. To compute grams of alcohol per day: (1) for each source of alcohol, multiply by the given factor and divide by 7 (if input is weekly intake) or 30 (if input is monthly intake) to get units/day; (2) multiply units/day by 8 to obtain grams/day; (3) sum grams/day intake of each source of alcohol to get total grams of alcohol per day.

Source	Factor
Red wine intake	1.5
Champagne/white wine	1.5
Beer/cider	2.5
Spirits	1
Fortified wine	1
Other alcoholic drinks	1.5

Table S4. The number of SNPs included in each PRS after removing SNPs in linkage disequilibrium via clumping.

Trait	# SNPs
Alcohol consumption	58
Alzheimer's disease	31
BMI	1,458
Breast cancer	153
CAD	207
Chronic kidney disease	4
Cirrhosis	2
Colorectal cancer	34
COPD	20
DBP	352
eGFR	31
Fasting blood glucose	24
HDL cholesterol	223
Hypertension	7
LDL cholesterol	195
Lung cancer	17
Pancreatic cancer	18
Parkinson's disease	44
Prostate cancer	123
SBP	390
Sleep duration	95
Smoking status	127
Stroke	79
Total cholesterol	240
Triglycerides	138
Type 2 diabetes	175
Total number of unique SNPs	3,941

BMI: body mass index; CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; HDL: high-density lipoprotein; LDL: low-density lipoprotein; SBP: systolic blood pressure; SNP: single nucleotide polymorphism; PRS: polygenic risk score.

Table S5. Summary statistics for the full cohort. Individuals who were related, were not of British ancestry, or had withdrawn their consent to participate were removed.

	Women	Men
Deaths by cause (n)		
Alzheimer's disease	43	38
with prevalent disease	0	1
without prevalent disease	43	37
Breast cancer	609	0
with prevalent disease	384	0
without prevalent disease	225	0
CAD	264	1,267
with prevalent disease	68	486
without prevalent disease	196	781
Chronic kidney disease	2	6
with prevalent disease	1	1
without prevalent disease	1	5
Cirrhosis	4	21
with prevalent disease	0	2
without prevalent disease	4	19
Colorectal cancer	295	445
with prevalent disease	35	94
without prevalent disease	260	351
COPD	119	218
with prevalent disease	74	126
without prevalent disease	45	92
Hypertension	4	9
with prevalent disease	4	9
without prevalent disease	0	0
Lung cancer	592	753
with prevalent disease	26	31
without prevalent disease	566	722
Pancreatic cancer	249	301
with prevalent disease	4	12
without prevalent disease	245	289
Parkinson's disease	18	64
with prevalent disease	9	41
without prevalent disease	9	23
Prostate cancer	0	436
with prevalent disease	0	183
without prevalent disease	0	253
Stroke	199	229
with prevalent disease	13	32
without prevalent disease	186	197
Type 2 diabetes	10	19
with prevalent disease	4	10
without prevalent disease	6	9
Prevalent disease (n)		
Alzheimer's disease	4	7
Breast cancer	6,323	0
CAD	5,445	12,530
Chronic kidney disease	170	311
Cirrhosis	27	70
Colorectal cancer	736	1,009
COPD	3,115	3,450
Hypertension	85,464	95,002
Lung cancer	107	122
Pancreatic cancer	15	26

Parkinson's disease	230	405
Prostate cancer	0	2,382
Stroke	2,126	3,092
Type 2 diabetes	4,072	7,576
Incident disease (n)		
Alzheimer's disease	314	345
Breast cancer	4,082	0
CAD	4,966	9,070
Chronic kidney disease	2,512	2,912
Cirrhosis	48	235
Colorectal cancer	1,036	1,437
COPD	2,740	3,576
Hypertension	3,154	3,371
Lung cancer	790	907
Pancreatic cancer	230	266
Parkinson's disease	299	493
Prostate cancer	0	4,542
Stroke	1,491	2,140
Type 2 diabetes	3,080	4,392
Mortality risk factors (mean (SD))		
Alcohol consumption (grams/day)	13.55 (12.34)	27.19 (23.39)
BMI (kg/m ²)	27.03 (5.14)	27.82 (4.21)
DBP (mmHg)	80.63 (9.93)	84.14 (9.99)
eGFR (mL/min/1.73 m ²)	85.57 (16.23)	87.61 (16.63)
Blood glucose (mmol/L)	5.07 (1.04)	5.18 (1.37)
HDL cholesterol (mmol/L)	1.60 (0.38)	1.28 (0.31)
LDL cholesterol (mmol/L)	3.64 (0.87)	3.49 (0.86)
SBP (mmHg)	135.60 (19.21)	141.30 (17.44)
Sleep duration (hours/day)	7.19 (1.10)	7.15 (1.07)
Smoking status (# ever smokers (%))	73,159 (40.5%)	79,226 (50.9%)
Total cholesterol (mmol/L)	5.90 (1.13)	5.50 (1.13)
Triglycerides (mmol/L)	1.56 (0.86)	1.98 (1.14)

CAD: coronary artery disease; COPD: chronic obstructive pulmonary disease; BMI: body mass index; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; HDL: high-density lipoprotein; LDL: low-density lipoprotein; SBP: systolic blood pressure; SD: standard deviation

Table S6. The estimated association between each mortality risk factor PRS and the risk factor measured at study baseline in women and men. Estimates are based on sex-specific linear regression models with robust standard error estimates, with the exception of smoking status, which was modeled using sex-specific logistic regression models. All models included adjustment for age at entry. Estimates are reported per standard deviation of the PRS.

Mortality risk factor	Women	Men
Alcohol consumption (grams/day)	0.90 (0.83, 0.96)	1.90 (1.79, 2.02)
BMI (kg/m ²)	1.46 (1.44, 1.49)	1.26 (1.24, 1.28)
DBP (mm Hg)	1.90 (1.85, 1.95)	1.63 (1.58, 1.68)
eGFR (mL/min/1.73 m ²)	2.54 (2.47, 2.61)	2.36 (2.28, 2.44)
Blood glucose (mmol/L)	0.065 (0.060, 0.070)	0.077 (0.069, 0.084)
HDL cholesterol (mmol/L)	0.118 (0.117, 0.120)	0.095 (0.094, 0.097)
LDL cholesterol (mmol/L)	0.234 (0.230, 0.238)	0.194 (0.190, 0.198)
SBP (mm Hg)	3.82 (3.74, 3.90)	3.06 (2.98, 3.14)
Sleep duration (hour)	0.092 (0.087, 0.097)	0.082 (0.077, 0.087)
Smoking status (odds ratio for ever smoking)	1.20 (1.19, 1.22)	1.22 (1.21, 1.23)
Total cholesterol (mmol/L)	0.300 (0.295, 0.305)	0.257 (0.251, 0.262)
Triglycerides (mmol/L)	0.187 (0.183, 0.191)	0.269 (0.263, 0.275)

BMI: body mass index; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; HDL: high-density lipoprotein; LDL: low-density lipoprotein; SBP: systolic blood pressure; PRS: polygenic risk score.

Table S7. The results of the analysis of all-cause mortality and the cPRS fitted in the training data and evaluated in the healthy subset of the test data. The cPRS were evaluated by fitting sex-specific Cox proportional hazards models of the association between age at death from all causes and the cPRS in the healthy subset of the test data. The healthy subset of the test data was defined as the test data with individuals with any of the diseases included as a top cause of death at baseline (prevalent cases). Both the continuous cPRS and categorical cPRS were modeled. The estimated HRs and CIs were converted to estimated years of life lost.

	Women	Men
Population in test data: N (deaths)		
Total population	29,379 (444)	18,249 (531)
Top 5% of cPRS	1,371 (27)	588 (21)
Middle 20% of cPRS	5,843 (80)	3,680 (107)
Bottom 5% of cPRS	1,647 (21)	1,145 (28)
Summary statistics for test data		
Age at entry (years; mean (SD))	54.4 (7.9)	54.3 (8.2)
Follow-up (years; mean (SD))	8.9 (0.9)	8.8 (1.1)
cPRS: HR (95% CI)		
Per SD of cPRS	1.07 (0.98, 1.18)	1.15 (1.06, 1.26)
Top 5% vs. middle 20% of cPRS	1.46 (0.94, 2.25)	1.28 (0.80, 2.04)
Bottom 5% vs. middle 20% of cPRS	0.89 (0.55, 1.44)	0.78 (0.51, 1.18)
cPRS: years of life lost (95% CI)		
Per SD of cPRS	0.71 (-0.21, 1.63)	1.43 (0.56, 2.31)
Top 5% vs. middle 20% of cPRS	3.75 (-0.61, 8.12)	2.45 (-2.23, 7.13)
Bottom 5% vs. middle 20% of cPRS	-1.16 (-5.97, 3.64)	-2.50 (-6.67, 1.66)

HR: hazard ratio; CI: confidence interval; cPRS: composite PRS; PRS: polygenic risk score; SD: standard deviation.

Table S8. Comparison of more sophisticated PRS with PRS based on GWAS Catalog. We compared the performance of three approaches for constructing PRS: (a) using results from the GWAS Catalog, (b) LD clumping and thresholding (C+T), and (c) LDpred. For the binary traits, we present the AUC for incident disease, while for the continuous risk factors, we present the R² for the measurement at baseline. See the Supplemental Methods for more details.

Trait (Reference) ^{a,b}	Best AUC or R ² from C+T/LDpred	GWAS Catalog AUC or R ²
Alzheimer's disease (2)	0.67	0.66
Breast cancer (3)	0.64	0.62
Chronic kidney disease (4)	0.54	0.53
Coronary artery disease (5)	0.57	0.58
Prostate cancer (6)	0.63	0.65
Stroke (7)	0.55	0.52 ^c
Type 2 diabetes (8)	0.61	0.61
BMI (9)	0.063	0.083
HDL cholesterol (10)	0.004	0.079
LDL cholesterol (10)	0.015	0.061
Total cholesterol (10)	0.009	0.060
Triglycerides (10)	< 0.001	0.048
Fasting plasma glucose (11)	0.006	0.003
eGFR (12)	0.027	0.022

^aSummary statistics for several of the disease traits were not available.

^bAvailable summary statistics for several of the mortality risk factors were based on data from the UK Biobank, precluding evaluation of the performance of the PRS in our data.

^cPRS based on recently published stroke PRS, not results from the GWAS Catalog (13).

AUC: area under the receiver operating characteristic curve; BMI: body mass index; C+T: clumping and thresholding; COPD: chronic obstructive pulmonary disease; DBP: diastolic blood pressure; eGFR: estimated glomerular filtration rate; GWAS: genome-wide association study; HDL: high-density lipoprotein; LD: linkage disequilibrium; LDL: low-density lipoprotein; PRS: polygenic risk score; SBP: systolic blood pressure; SNP: single nucleotide polymorphism.

SUPPLEMENTAL METHODS

Constructing Alternative PRS

In addition to the PRS based on genome-wide significant results from the GWAS Catalog, which were the focus of the paper, we considered PRS constructed by two alternative approaches: linkage disequilibrium (LD) clumping and thresholding (C+T) (14) and LDpred (15). Briefly, C+T involves first clumping variants so as to remove correlated SNPs. This is accomplished by selecting the most significant variant and removing nearby variants that are correlated beyond some r^2 . Next, the clumped variants are thresholded, i.e., those with p-values larger than some p are removed. LDpred is a Bayesian method that leverages information on LD structure to infer the posterior mean effect size of each SNP. This method utilizes a prior effect size distribution which requires specification of the proportion of causal SNPs, q . Thus, C+T involves selecting two tuning parameters, r^2 and p , while LDpred requires selecting one tuning parameter, q .

We tuned the parameters for C+T and LDpred using data from the UK Biobank. In particular, we considered $r^2 = (0.1, 0.2, 0.4, 0.6, 0.8)$, $p = (5 \times 10^{-8}, 5 \times 10^{-6}, 5 \times 10^{-4}, 0.05, 1)$, and $q = (0, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1)$. For the binary disease traits, we chose the tuning parameters based on the AUC for prevalent disease. For the continuous risk factors, we chose the tuning parameters based on the R^2 for the baseline risk factor measurements.

Evaluating Alternative PRS

For each trait, we compared the three PRS (GWAS Catalog, C+T, and LDpred) in terms of their ability to predict the trait. Specifically, for the binary disease traits, we evaluated the AUC for incident disease in the UK Biobank. For the continuous risk factors, we evaluated the R^2 for the baseline risk factor measurements in the UK Biobank. We compared the best AUC/ R^2 among the PRS constructed by C+T and LDpred to the AUC/ R^2 for the GWAS Catalog PRS for the same trait.

Mortality Analysis with Alternative PRS

The performance of the GWAS Catalog PRS in terms of the AUC/ R^2 was generally quite good (relative to the PRS built using C+T or LDpred), though we found meaningful gains for the LDpred PRS for breast cancer and stroke (Table S8). Thus, we repeated the main composite PRS mortality analysis, using the LDpred PRS for breast cancer and stroke (with $q = 0.1$ and 0.03 , respectively) in place of the GWAS Catalog versions of these PRS. In particular, we constructed a new composite PRS in the training data using the LDpred PRS for breast cancer and stroke and the GWAS Catalog PRS for the other traits. All other aspects of the construction of this composite PRS were as described for the main composite PRS mortality analysis. We then evaluated the association between the new composite PRS and all-cause mortality by estimating the hazard ratio per standard deviation of the composite PRS in the test data, adjusting for the first ten principal components.

SUPPLEMENTAL REFERENCES

1. Eastwood S V., Mathur R, Atkinson M, Brophy S, Sudlow C, Flaig R, De Lusignan S, Allen N, Chaturvedi N. Algorithms for the capture and adjudication of prevalent and incident diabetes in UK Biobank. *PLoS One*. 2016;11(9).
2. Lambert JC, Ibrahim-Verbaas CA, Harold D, Naj AC, Sims R, Bellenguez C, Jun G, DeStefano AL, Bis JC, Beecham GW, et al. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. *Nat Genet*. 2013;45(12):1452–8.
3. Michailidou K, Lindström S, Dennis J, Beesley J, Hui S, Kar S, Lemaçon A, Soucy P, Glubb D, Rostamianfar A, et al. Association analysis identifies 65 new breast cancer risk loci. *Nature*. 2017;551(7678):92–4.
4. Pattaro C, Teumer A, Gorski M, Chu AY, Li M, Mijatovic V, Garnaas M, Tin A, Sorice R, Li Y, et al. Genetic associations at 53 loci highlight cell types and biological pathways relevant for kidney function. *Nat Commun*. 2016;7(1):1–19.
5. Nikpay M, Goel A, Won HH, Hall LM, Willenborg C, Kanoni S, Saleheen D, Kyriakou T, Nelson CP, CHopewell J, et al. A comprehensive 1000 Genomes-based genome-wide association meta-analysis of coronary artery disease. *Nat Genet*. 2015;47(10):1121–30.
6. Schumacher FR, Al Olama AA, Berndt SI, Benlloch S, Ahmed M, Saunders EJ, Dadaev T, Leongamornlert D, Anokian E, Cieza-Borrella C, et al. Association analyses of more than 140,000 men identify 63 new prostate cancer susceptibility loci. *Nat Genet*. 2018;50(7):928–36.
7. Malik R, Chauhan G, Traylor M, Sargurupremraj M, Okada Y, Mishra A, Rutten-Jacobs L, Giese AK, Van Der Laan SW, Gretarsdottir S, et al. Multiancestry genome-wide association study of 520,000 subjects identifies 32 loci associated with stroke and stroke subtypes. *Nat Genet*. 2018;50(4):524–37.
8. Scott RA, Scott LJ, Mägi R, Marullo L, Gaulton KJ, Kaakinen M, Pervjakova N, Pers TH, Johnson AD, Eicher JD, et al. An Expanded Genome-Wide Association Study of Type 2 Diabetes in Europeans. *Diabetes*. 2017;66(11):2888–902.
9. Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR, Powell C, Vedantam S, Buchkovich ML, Yang J, et al. Genetic studies of body mass index yield new insights for obesity biology. *Nature*. 2015;518(7538):197–206.
10. Willer CJ, Schmidt EM, Sengupta S, Peloso GM, Gustafsson S, Kanoni S, Ganna A, Chen J, Buchkovich ML, Mora S, et al. Discovery and refinement of loci associated with lipid levels. *Nat Genet*. 2013;45(11):1274–85.
11. Scott RA, Lagou V, Welch RP, Wheeler E, Montasser ME, Luan J, MäGi R, Strawbridge RJ, Rehnberg E, Gustafsson S, et al. Large-scale association analyses identify new loci influencing glycemic traits and provide insight into the underlying biological pathways. *Nat Genet*. 2012;44(9):991–1005.
12. Li M, Li Y, Weeks O, Mijatovic V, Teumer A, Huffman JE, Tromp G, Fuchsberger C, Gorski M, Lyytikäinen LP, et al. SOS2 and ACP1 loci identified through large-scale exome chip analysis regulate kidney development and function. *J Am Soc Nephrol*. 2017;28(3):981–94.
13. Rutten-Jacobs LCA, Larsson SC, Malik R, Rannikmäe K, Sudlow CL, Dichgans M, Markus HS, Traylor M. Genetic risk, incident stroke, and the benefits of adhering to a healthy lifestyle: Cohort study of 306 473 UK Biobank participants. *BMJ*. 2018;363.
14. Privé F, Vilhjálmsson BJ, Aschard H, Blum MGB. Making the most of clumping and thresholding for polygenic scores. *Am J Hum Genet*. 2019;105(6):1213–21.
15. Vilhjálmsson BJ, Yang J, Finucane HK, Gusev A, Lindström S, Ripke S, Genovese G, Loh PR, Bhatia G, Do R, et al. Modeling linkage disequilibrium increases accuracy of

polygenic risk scores. *Am J Hum Genet.* 2015;97(4):576–92.