

Full Title: A protocol for using human genetic data to identify circulating protein level changes that are the causal consequence of cancer processes.

Short title: Protocol for using human genetics to identify circulating proteins associated with cancer process.

Lisa M Hobson [1, 2]; Prof. Richard M Martin [2, 3]; Dr. Karl Smith-Byrne [3, 4]; Prof. George Davey Smith [1, 3]; Prof. Gibran Hemani [1, 3]; Dr. Joseph H Gilbody [1,2]; Dr. James Yarmolinsky [1, 2, 5]; Dr. Sarah ER Bailey [6]; Dr. Lucy J. Goudswaard [1, 2]; Dr. Philip C Haycock [1]

1. MRC Integrative Epidemiology Unit, Bristol Medical School, University of Bristol
2. Population Health Sciences, Bristol Medical School, University of Bristol
3. Bristol NIHR Biomedical Research Centre, University Hospitals Bristol and Weston NHS Foundation Trust and the University of Bristol
4. Cancer Epidemiology Unit, Nuffield Department of Population Health, University of Oxford
5. Department of Epidemiology and Biostatistics, School of Public Health, Imperial College London
6. Department of Health and Community Science, University of Exeter

Email Addresses

Lisa M Hobson: lisa.hobson@bristol.ac.uk

Richard M Martin: richard.martin@bristol.ac.uk

Karl Smith-Byrne: Karl.Smith-Byrne@ndph.ox.ac.uk

George Davey-Smith: KZ.Davey-Smith@bristol.ac.uk

Gibran Hemani: g.hemani@bristol.ac.uk

Joseph H Gilbody: joe.gilbody@bristol.ac.uk

James Yarmolinsky: j.yarmolinsky@imperial.ac.uk

Sarah ER Bailey: s.e.r.bailey@exeter.ac.uk

Lucy J Goudswaard: lucy.goudswaard@bristol.ac.uk

Philip C Haycock: philip.haycock@bristol.ac.uk

Corresponding Authors

Lisa M Hobson

MRC Integrative Epidemiology Unit, Oakfield House, Oakfield Grove, Clifton, Bristol BS8 2BN

lisa.hobson@bristol.ac.uk

Lucy J Goudswaard

MRC Integrative Epidemiology Unit, Oakfield House, Oakfield Grove, Clifton, Bristol BS8 2BN

lucy.goudswaard@bristol.ac.uk

Philip C Haycock

MRC Integrative Epidemiology Unit, Oakfield House, Oakfield Grove, Clifton, Bristol BS8 2BN

philip.haycock@bristol.ac.uk

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

Abstract

Introduction

Cancer is a leading cause of death worldwide. Early detection of cancer improves treatment options and patient survival but detecting cancer at the earliest stage presents challenges. Identification of circulating protein biomarkers for cancer risk stratification and early detection is an attractive avenue for potentially minimally invasive screening and early detection methods. We hypothesise that protein level changes resulting from cancer development can be identified via an individual's polygenic risk score (PRS) for the disease, representing their genetic liability to developing that cancer.

Methods and analysis

PRS will be calculated using the PRS continuous shrinkage approach (PRS-CS and PRS-CSx) for colorectal and lung cancer risk. This methodology utilises effect sizes from summary statistics from genome-wide association studies (GWAS) available for the cancers of interest to generate weights via the continuous shrinkage approach which incorporates the strengths of the GWAS associations into the shrinkage applied (1). This methodology both improves upon previous PRS methods in accuracy as well as improving cross-ancestry application in the PRS-CSx approach. GWAS summary statistics will be from the Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO) and the International Lung Cancer Consortium (ILCCO). The association between the polygenic risk scores and 2923 proteins measured by the Olink platform in UK Biobank (UKB) participants with protein measures available will be assessed using linear regression under the assumption of linearity in the proteomic data. The proteins identified could represent several different scenarios of association such as forward causation (protein causes cancer), reverse causation (cancer genetic liability causes protein level change), or horizontal pleiotropy bias (no causal relationship exists between the protein and cancer). Forward and reverse Mendelian randomization sensitivity analyses, as well as colocalization analysis, will be performed in efforts to distinguish between these three scenarios. Protein changes identified as causally downstream of genetic liability to cancer could reflect processes occurring prior to, or after, disease onset. Due to individuals in the UKB having proteins measures at only one timepoint, and because UKB contains a mix of incident and prevalent cases, some protein measures will have been made prior to a cancer diagnosis while others will have been made after a cancer diagnosis. We will explore the strength of association in relation to the time between protein measurement and prevalent or incident cancer diagnosis.

Ethics and Disseminations

No additional ethical approval is required for Genome Wide Association (GWAS) data used in this analysis as all data from GWAS has undergone individual ethical approval prior to this study. UK Biobank protein measure data will be obtained under application ID: 15825/81499.

Results produced from these analyses will be submitted as an open-access manuscript to journals for review and all code will be made publicly available using GitHub. The PRS we generate and the results of the PRS-protein associations will be returned to the UK Biobank.

Strengths and limitations of this study

- A strength of the proposed PRS method in this study is the use of all available SNPs from a GWAS, which may increase power to identify proteins in comparison with conventional Mendelian Randomisation (MR) methods that use only those SNPs that are genome-wide significant.
- Limitations of the study:
 - Lack of protein data for diverse population groups within available datasets; therefore, results may not be generalisable to ancestries outside of the European population for whom sufficient protein data was available for this study.
 - UKB participants reflect a subset of the population from a higher socioeconomic position than average.

- Prevalent cancer cases will reflect a specific subset of the general population with cancer, individuals who have survived cancer and were able to volunteer for the study; potentially introducing survivorship bias.
- It cannot be ruled out that proteins may reflect effects of processes beyond cancer liability to protein pathways.
- Lack of staging information for cancer cases within the UKB limiting our ability to distinguish early versus more advanced cancers.
- The proteomic technology currently used measures protein binding as opposed to protein levels

Introduction

Detecting cancer at an early stage is important because patients diagnosed early have a greater chance of being treated with curative intent and so experience increased long-term survival. Cancer is a leading cause of death worldwide (2) and 5-year survival rates fall considerably when diagnosis is made at a later stage. The 5-year survival of colorectal cancer and lung cancer is reduced from more than 9 in 10 and 6 in 10, respectively, when diagnosed at stage 1, to 1 in 10 for colorectal cancer and less than 1 in 10 for lung cancer when diagnosed at stage 4 (3). The NHS Long Term Plan aims to increase early detection of cancers from half to three quarters by the year 2028 to improve cancer survival (4) as currently in England only 54% of cancers are detected early (5). To meet this goal a number of research challenges need to be addressed, including development of methods for determining cancer risk (i.e. risk stratification) and identifying biomarkers which are effective for detecting cancer at an early stage (6).

One of the challenges to be overcome in improving cancer early detection is the identification of specific biomarkers for the cancers of interest that can be measured by minimally invasive, low-cost methods and are able to be implemented in a clinical setting. One way to address this challenge is the measurement of circulating protein levels in blood serum or plasma, potentially feasible because of the widespread use of blood tests in healthcare. Circulating protein biomarkers are a potentially useful tool for several clinical areas including identifying groups at high risk of the future development of cancer (risk stratification), early detection, disease diagnosis and monitoring biological processes (7,8). They may provide a minimally invasive means of screening asymptomatic individuals for undiagnosed disease or for diagnosis of symptomatic patients (9,10). Advantages of measuring protein within the blood, as opposed to other minimally invasive methods, such as circulating-tumour DNA (ctDNA) sequencing, is the reduced volume of sample required for analysis and affordability. For the Olink explore 3072 panel, only 6 μ L of plasma or serum is required vs. 4-5mL of plasma to obtain 5-10ng/mL of ctDNA and with the possibility of implementation of protein testing via ELISA, costing around ~£4 per test (11–16).

One approach to biomarker discovery is via prospective cohort studies to identify proteins associated with the incidence of a disease of interest, by measuring protein levels in individuals before diagnosis. These methods require large sample sizes over long periods of time to capture these events, at great financial and time cost (17). A comparatively inexpensive technique for biomarker discovery has been formalised by Holmes and Davey Smith (18), and involves application of Mendelian randomization (MR) of disease liability as the exposure on protein levels as the outcome (sometimes described as reverse MR or reverse gear MR). Building on this idea, we propose that protein level changes resulting from cancer onset can be identified via an individual's PRS for specific cancers, representing their genetic liability to developing that cancer. Defining the point of "cancer onset" remains difficult, with many possible mechanisms of initiation; for the purpose of this study we will use date of diagnosis to determine prevalent vs. incident cases within the cohort (19).

Proteins associated with genetic liability to cancer could reflect different mechanisms of association. Associations could reflect 'forward causation' where the protein is upstream of and causal for cancer or 'reverse causation' where carcinogenesis is causing the change in downstream protein level (Figure 1).

Proteins that are associated via forward causation are upstream of the cancer pathway and therefore do not always denote the presence of cancer but could identify potential therapeutic targets for disease prevention and cancer prediction, these protein levels will likely remain stable over long periods of time. Proteins downstream of cancer development will likely show more variation in levels as a result of the progression of cancer; we will refer to these proteins as “reverse causal”. For proteins that cause cancer, most of the variants in the cancer PRS will have no causal relationship to those proteins, whereas for proteins that are causally associated downstream of cancer liability pathways, all variants in the PRS could contribute to the association signal (Figure 1). We thus expect a cancer PRS to be better powered for discovery of proteins downstream of cancer development. However, the relative balance in findings reflecting scenarios one (‘forward causation’) and two (‘reverse causation’) is likely to depend on the prevalence of disease, including early or pre-clinical stages, in the sample used to measure the proteins. In general, the higher the prevalence, the greater the number of associations we expect to see reflecting effects of cancer liability pathways on protein concentration. Association of proteins with a genetic liability to cancer can also be due to factors other than genetic liability such as horizontal pleiotropy bias, as illustrated in Figure 1 by G4 and G6; which may negate its use in risk stratification or early detection.

Aims

The overall aim of these analyses is to identify protein changes that are the causal consequence of genetic liability to cancer.

Methods and Analysis

Polygenic Risk Score Analyses

PRS can be developed using GWAS summary statistics on the associations of many SNPs across the genome with cancer. In this way millions of SNPs can be combined to develop an individual’s PRS for a disease. A PRS is the sum of the number of copies of risk alleles individuals have for SNPs across the genome, weighted by the effect size of these SNPs in relation to the disease of interest (20). While the initiation of cancer and the factors that contribute to the onset and progression of cancer are still not fully understood (21); by calculating a PRS using data from all SNPs across the genome, SNPs involved in initiation, promotion and progression of cancer will be captured by this score, reflecting the complex process of cancer development (19).

In this study, PRS will be calculated for UKB participants (Application ID: 15825/81499) with proteomic measurements (N=49,542), individuals with sex-mismatch (derived by comparing genetic sex and reported sex) or individuals with sex-chromosome aneuploidy will be excluded from the analysis (N=814) as well as highly related individuals related to a 3rd degree to >200 individuals (N=2) (22). For colorectal cancer, we will use effect weights derived from GWAS summary statistics of the: i) Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO) (GWAS Catalog Accession: GCST90255675) for Europeans; and ii) the Asia Colorectal Cancer Consortium (ACCC)/Korean-National Cancer Centre CRC Study 2 (Korea-NCC2) for GWAS summary statistics for East Asians. For lung cancer, we will use effect weights derived from GWAS summary statistics from the International Lung Cancer Consortium (ILCCO) (GWAS Catalog Accession: GCST004748). Sample selection and quality control within these studies has been previously described (23,24).

PRS will be derived from the PRS-CS and PRS-CSx approaches, using summary statistics from GWAS for the cancer of interest along with an external linkage disequilibrium (LD) reference panel corresponding to the ancestry of the GWAS. The continuous shrinkage approach incorporates the strengths of the GWAS associations into the shrinkage applied to shrink small SNP effects towards zero, while large effects are unaffected (25), generating a posterior effect size for each SNP (1). These weights will be used to calculate the PRS of UK Biobank participants for colorectal and lung cancer, calculating the sum of risk increasing alleles across all genetic variants weighted by the effect sizes generated by PRS-CS (1,26). PRS-CSx applies the same methodology as PRS-CS to multi-ancestry GWAS summary statistics, improving generalisability of results to more ancestry groups within the global majority. In an effort to reduce the Eurocentric bias and

to increase power, we will be utilising GWAS summary statistics for colorectal cancer from European and East Asian ancestries to develop polygenic risk scores (23,27,28).

Cancer Subgroup Analyses

In addition to a PRS for overall colorectal and lung cancers, we will calculate a PRS for colon cancer and rectal cancer specifically and for lung cancer subgroups (adenocarcinoma, squamous cell carcinoma, small cell carcinoma). Additionally, PRS scores will be calculated for never smokers and ever smokers using weights generated from summary statistics of GWAS for lung cancer in never smokers and lung cancer in ever smokers.

Olink Proteins

Olink protein measurements were performed as part of the Pharma Proteomic Project (UKB-PPP) on blood plasma samples using the antibody-based protein Olink Explore 3072 Proximity Extension Assay. Proteomics were generated for 54,219 participants considered to be highly representative of the UK Biobank population on baseline characteristics, enriched for selected diseases (29). The number of participants with colorectal and lung cancers can be seen in Table 1. Quality control, sample selection and data processing has been described previously (29). Associations between the participants' PRS and 2923 Olink protein measures from the UK Biobank will be tested via linear regression, adjusting for age, sex, principle components and sample storage time where this has an impact on protein level variation (30). Protein measures will undergo inverse rank normal transformation (INT) for each protein (31). The number of independent proteins will be calculated using the metaboprep R package (32). False discovery rate correction will be applied to p-values, proteins with p-value less than the calculated alpha will be prioritised for further analyses.

Table 1. Disease frequency within the UKB cohort and within the UKB-PPP study participants.

ICD10/ICD9 code	Disease	Number of cases (UKB)	Number of cases (UKB-PPP)
C18/153	Malignant neoplasm of colon	5751 (1.14%)	602 (1.13%)
C19/1540	Malignant neoplasm of rectosigmoid junction	635 (0.13%)	65 (0.12%)
C20/1541	Malignant neoplasm of rectum	2562 (0.51%)	254 (0.48%)
C34/162	Malignant neoplasm of bronchus and lung	4917 (0.98%)	552 (1.04%)
ICD-O-3 Code	Histological Subset	Number of cases (UKB)	Number of cases (UKB-PPP)
8140, 8211, 8250–8260, 8310, 8323, 8480–8490, 8550	Lung adenocarcinoma	2200 (0.43%)	238 (0.45%)
8070-8072	Lung squamous cell carcinoma	1164 (0.23%)	142 (0.27%)
8041–8042	Lung small cell carcinoma	465 (0.09%)	60 (0.11%)

Number of cases for each cancer type derived from UK biobank phenotypic data with percentage of cases out of overall individuals in represented in brackets. UKB overall n = 502,178, UKB with protein measures (UKB-PPP) n = 53,058.

Table 2. Prevalent and incident cases from UKB cohort and within the UKB-PPP study participants.

ICD10/ICD9 Code	Disease	UKB Overall	UKB-PPP
-----------------	---------	-------------	---------

		Prevalent Cases	Incident Cases	Prevalent Cases	Incident Cases
C18/153	Malignant neoplasm of colon	1098	4653	105	497
C19/1540	Malignant neoplasm of rectosigmoid junction	441	1964	23	42
C20/1541	Malignant neoplasm of rectum	598	1964	63	191
C34/162	Malignant neoplasm of bronchus and lung	248	4669	31	521
ICD-O-3 Code	Histological Subset	Prevalent Cases	Incident Cases	Prevalent Cases	Incident Cases
8140, 8211, 8250–8260, 8310, 8323, 8480–8490, 8550	Lung adenocarcinoma	78	2122	8	230
8070-8072	Lung squamous cell carcinoma	75	1089	14	128
8041–8042	Lung small cell carcinoma	24	441	3	57

Number of incident and prevalent cases for each cancer type derived from UK biobank phenotypic data.

Sensitivity Analyses

Proteins identified from association analyses may reflect different scenarios, including causation or confounding from population stratification or dynastic effects. Some possible scenarios include: (1) a protein may be a cause of cancer risk, which we define as “forward causation”; (2) an alternative scenario is that the protein identified is causally downstream of cancer liability, which we refer to as “reverse causation”; (3) there is no causal relationship between the protein and cancer and the identified association reflects horizontal pleiotropy, (4) due to population stratification where spurious associations are due to differences in the GWAS population and those that the PRS is calculated on. We will perform various sensitivity analyses to distinguish amongst these scenarios, described below.

Bidirectional Mendelian Randomisation Sensitivity Analyses

MR uses genetic variants, associated with the phenotype of interest as the instrumental variable to assess the effect of the phenotype on an outcome. Due to the random nature of inheritance of genetic variants there is an advantage over observational epidemiology whereby confounders may influence both the exposure and outcome of interest (33–35). Genetic associations used in MR analyses often come from GWAS summary data, whereby association is conventionally defined by a p-value threshold of 5×10^{-8} .

Assumptions

The three core assumptions of MR, known as the instrumental variable (IV) assumptions (Figure 2), are relevance (IV1) – is the instrumental variable (G) associated with the exposure (E), independence (IV2) – there is no confounding of the association between the instrument (G) and outcome (O) (this can arise through population stratification, dynastic effects and assortative mating) and exclusion restriction (IV3) – the instrumental variable (G) does not act on the outcome (O) except via the exposure (E) e.g. no horizontal pleiotropy (red dashed line, IV3); (36–38).

Study Design

MR will be performed in the forward and reverse direction: forward MR will be used to estimate the effect of selected proteins on the cancer of interest and reverse MR will be used to estimate the effect of cancer liability on circulating protein concentration (18). Forward MR will be performed using cis-pQTLs to instrument proteins identified as being associated with the cancer PRS, the threshold for these will be $p < 3.4 \times 10^{-11}$ (39). Cis-pQTLs will be defined as within < 1 Mbp of the protein coding gene and trans-pQTLs will be defined as > 1 Mbp away from the protein coding gene (40). Reverse MR will be performed using SNPs associated with the cancer PRS at a threshold of $p < 5 \times 10^{-8}$. If association is found in the forward direction this may suggest that the protein is causal for the disease but if association is found in the reverse direction

this may suggest that genetic liability to cancer is causing the protein level change (18); to elucidate this causality, different MR estimation methods will be employed, the application and conditions of these are described below.

Instrument & Method Selection

The strength of instrument will be determined by calculating the F-statistic, a measure of potential weak instrument bias that could arise from the use of IVs as a proxy for the effect of exposure on outcome (38). The F-statistic takes into account the genetic variance (R^2), the sample size and how many instruments are present. An F-statistic greater than 10 indicates that the bias from weak instruments is small, where this F-stat is less than 10 this indicates a possibility of bias and will be noted (41).

Dependent on the number of SNPs available, the appropriate method of effect-estimation will be selected for MR analyses. For proteins with a single pQTL SNP the Wald ratio will be calculated as the ratio of SNP-outcome/SNP-exposure association (42). For proteins with two or three independent SNPs, a fixed-effects inverse variance weighted (IVW) model will be used. For four or more independent SNPs, a random effects inverse variance (IVW) model, combining multiple SNP outcome/exposure Wald ratio, will be used (43). In the event of multiple independent SNPs, pleiotropy will be considered by calculating Cochran's Q statistic, a method for assessing global and individual pleiotropy across instruments (44). Weighted mode and weighted median methods will also be used when > 10 SNPs are available (45,46).

The MR-PRESSO, weighted mode and weighted median methods will be used to assess IVs for horizontal pleiotropy, violation of IV3 where the IV acts on the outcome not via the exposure, by comparing estimates with and without suspected pleiotropic variants, this will be repeated for both forward and reverse MR (45,46, 47). In addition being robust to pleiotropy methods such as MR robust adjusted profile score (RAPS) also accounts for other potential sources of bias such as weak instruments and measurement error in the exposure (48). MR-CAUSE (Causal Analysis Using Summary Effect estimates) is another method that can be used when IV3 is violated due to pleiotropic effects of correlated pleiotropy, where the pleiotropic factor is a confounder of the exposure-outcome association versus when the IV has effects on pleiotropic factor independently of the effect of the IV on the exposure this is uncorrelated pleiotropy (49).

When performing reverse MR using a larger numbers of SNPs, clustered heterogeneity can occur when different genetic variants are causally associated via distinct pathways, to assess this, clustering based methods can be used to divide groups based on these estimates of causality (50). MR-Clust will be used to investigate clustered heterogeneity across IVs and identify potential distinct pathways that make up the effect estimate; clustering works by separating the variants into clusters with additional null and junk clusters, representing no causal effect or those that do not fit within the distinct clusters (50,51). Another clustering method that will be used is the Noise-Augmented von Mises-Fisher Mixture model (NAvMix), this method allows for variants to belong to multiple clusters based on their probability of membership to that cluster (52,53). The contamination mixture model method can also be used to cluster into distinct groups based on the IVs causal effect estimate even when invalid IVs are present (54). PheWAS-based clustering will also be used to cluster SNP associations based on different pathways and thus help identify other causal pathways of the PRS – protein associations found (55). The methods of MR described make different assumptions and aim to address different violation of the IV assumptions, testing of these different methods has illustrated the variation in accuracy and the need for appropriate method selection based on the datasets used (56).

Data Harmonisation

Harmonisation of the effect alleles across the GWAS summary statistics datasets will be performed using the TwoSampleMR R package (57).

Colocalization

Associations may be due to genomic confounding, where genetic variants in linkage disequilibrium (LD) at the same locus act on the cancer and protein via separate pathways, a form of horizontal pleiotropy bias. Colocalization analyses will be used to assess if genetic associations with cancer and proteins are due to shared causal variants at the same locus through genomic confounding (58,59).

Further Analyses

“Time-to” and “Time-from” Diagnosis

In observational analyses, we will evaluate the magnitude of the relationship between proteins, taken either pre or post-diagnosis, and cancer risk. This will involve an analysis of prevalent and incident cancer cases in UKB (Table 2) and a time variable derived from date of cancer diagnosis and time of blood collection (60,61). To adjust for any variation in protein concentration as a result of sample storage and protein degradation over time (30), the relationship between storage time and protein level for all protein measures available will be assessed. Proteins are more likely to be causally downstream of cancer onset if the association with cancer is sensitive to time between protein measure and cancer diagnosis, a potential route for differentiating between normal baseline levels and levels that suggest the presence of cancer. If protein levels are detectable prior to patient reported symptoms proteins may be more suited for screening and early detection.

Replication of Findings

Replication of protein association and MR will be carried out in the DECODE cohort (62) and EPIC study (63) where proteins are available.

Software

This work will be carried out using the computational facilities of the Advanced Computing Research Centre, University of Bristol - <http://www.bristol.ac.uk/acrc/>.

PRS-CS (<https://github.com/getian107/PRScs>) and PRS-CSx (<https://github.com/getian107/PRScsx>) will be used to calculate polygenic risk scores, using R, Python and PLINK.

Metaboprep (<https://github.com/MRCIEU/metaboprep>) R package will be used to calculate independent proteins (32). Mendelian Randomization analyses and data harmonisation will be performed using the R packages TwoSampleMR (<https://github.com/MRCIEU/TwoSampleMR>) and MendelianRandomization (<https://cran.r-project.org/web/packages/MendelianRandomization/index.html>) (64).

Proteins will be inverse rank normal transformed using the “RankNorm” function in R package “RNOmni” (<https://cran.r-project.org/web/packages/RNOmni/index.html>) (31).

Patient and Public Involvement

A summary of the proposed research was presented to members of a patient and public involvement group, with either personal experience with cancer or experience via a family member. The feedback received was that this was very important research and that they believe it would be useful for early detection for cancers that do not yet have specific screening via a blood test. Updates about this study will also be disseminated to the group.

Ethics and Dissemination

The colorectal cancer GWAS conducted by Fernandez-Rozadilla et al. (2022) was approved by the South Central Ethics Committee (UK) under the reference number 17/SC/0079 (23).

All studies used in the lung cancer GWAS conducted by McKay et al. (2017) obtained local ethics committee approval and all participants gave informed consent (24).

Application for colorectal cancer site specific GWAS summary statistics from GECCO has been submitted.

Application for summary statistics from the Asian Colorectal Cancer Consortium (ACCC) and the Korean-National Cancer Center CRC Study 2 (Korea-NCC2) will be submitted.

UK Biobank was approved by the North West Multi-centre Research Ethics Committee (MREC) as a Research Tissue Bank (RTB) approval renewed in 2021, all participants in the study have given informed consent (65). Genotype, phenotype and Olink protein measure data access has been obtained under Application ID: 15825/81499.

Results of these analyses will be disseminated via the University of Bristol MRC Integrative Epidemiology Unit IEU Portal and submitted as a manuscript to a peer-reviewed journal for publication. All statistical code will be made available via GitHub.

Polygenic risk scores and PRS-protein associations will be returned to the UK Biobank in line with the UK Biobank obligation for researchers outlined (66).

Data Availability Statement

No data has been collected or generated as part of this protocol.

Funding Statement

LMH is supported in part by grant MR/N0137941/1 for the GW4 BIOMED MRC DTP, awarded to the Universities of Bath, Bristol, Cardiff and Exeter from the Medical Research Council (MRC)/UKRI RMM is a National Institute for Health Research Senior Investigator (NIHR202411). RMM, LJG and PCH are supported by a Cancer Research UK 25 (C18281/A29019) programme grant (the Integrative Cancer Epidemiology Programme). RMM is also supported by the NIHR Bristol Biomedical Research Centre which is funded by the NIHR (BRC-1215-20011) and is a partnership between University Hospitals Bristol and Weston NHS Foundation Trust and the University of Bristol. Department of Health and Social Care disclaimer: The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care.

SERB was supported by an NIHR Advanced Fellowship (NIHR 301666) whilst undertaking this work. Additional support was provided by the Higgins family.

Competing Interests

RMM, LJG and PCH have received funding from Cancer Research UK. LMH receives funding from the GW4 BioMed2 MRC DTP.

Author Contributions

LMH: Writing – Original draft

PCH, LJG, KSB, RMM, SERB: Conceptualisation, Methodology, Writing – Review & Editing, Supervision

GH: Writing – Review & Editing, Diagram

GDS, JHG & JY: Writing – Review & Editing

References

1. Ge T, Chen CY, Ni Y, Feng YCA, Smoller JW. Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nat Commun* [Internet]. 2019 Apr 16 [cited 2023 Oct 5];10(1):1776. Available from: <https://www.nature.com/articles/s41467-019-09718-5>
2. WHO. Cancer [Internet]. 2022 [cited 2023 Nov 1]. Available from: <https://www.who.int/news-room/fact-sheets/detail/cancer>

3. Cancer Research UK [Internet]. 2015 [cited 2023 Dec 9]. Why is early cancer diagnosis important? Available from: <https://www.cancerresearchuk.org/https%3A//www.cancerresearchuk.org/about-cancer/spot-cancer-early/why-is-early-diagnosis-important>
4. NHS England. NHS Long Term Plan [Internet]. NHS England; 2019 Jan [cited 2023 Nov 7]. Report No.: 1.2. Available from: <https://www.longtermplan.nhs.uk/publication/nhs-long-term-plan/>
5. Health Education England. Health Education England. 2023 [cited 2023 Dec 9]. Improving cancer diagnosis and earlier detection. Available from: <https://www.hee.nhs.uk/our-work/primary-care/improving-cancer-diagnosis-earlier-detection>
6. Crosby D, Bhatia S, Brindle KM, Coussens LM, Dive C, Emberton M, et al. Early detection of cancer. *Science* [Internet]. 2022 Mar 18 [cited 2023 Oct 5];375(6586):eaay9040. Available from: <https://www.science.org/doi/10.1126/science.aay9040>
7. Integrative Analysis of Lung Cancer Etiology and Risk (INTEGRAL) Consortium for Early Detection of Lung Cancer, Guida F, Sun N, Bantis LE, Muller DC, Li P, et al. Assessment of Lung Cancer Risk on the Basis of a Biomarker Panel of Circulating Proteins. *JAMA Oncol*. 2018 Oct 1;4(10):e182078.
8. Carrasco-Zanini J, Pietzner M, Davitte J, Surendran P, Croteau-Chonka DC, Robins C, et al. Proteomic prediction of common and rare diseases [Internet]. medRxiv; 2023 [cited 2023 Dec 6]. p. 2023.07.18.23292811. Available from: <https://www.medrxiv.org/content/10.1101/2023.07.18.23292811v1>
9. Pavlou MP, Diamandis EP. The cancer cell secretome: A good source for discovering biomarkers? *J Proteomics* [Internet]. 2010 Sep 10 [cited 2023 Nov 21];73(10):1896–906. Available from: <https://www.sciencedirect.com/science/article/pii/S1874391910001296>
10. Califf RM. Biomarker definitions and their applications. *Exp Biol Med* [Internet]. 2018 Feb [cited 2023 Nov 21];243(3):213–21. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5813875/>
11. Gao Q, Zeng Q, Wang Z, Li C, Xu Y, Cui P, et al. Circulating cell-free DNA for cancer early detection. *The Innovation* [Internet]. 2022 May 6 [cited 2024 Apr 5];3(4):100259. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9133648/>
12. Yan Y yan, Guo Q ru, Wang F hua, Adhikari R, Zhu Z yan, Zhang H yan, et al. Cell-Free DNA: Hope and Potential Application in Cancer. *Front Cell Dev Biol* [Internet]. 2021 Feb 22 [cited 2024 Apr 5];9:639233. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7938321/>
13. Song P, Wu LR, Yan YH, Zhang JX, Chu T, Kwong LN, et al. Limitations and opportunities of technologies for the analysis of cell-free DNA in cancer diagnostics. *Nat Biomed Eng* [Internet]. 2022 Mar [cited 2024 Apr 5];6(3):232–45. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9336539/>
14. Kim H, Park KU. Clinical Circulating Tumor DNA Testing for Precision Oncology. *Cancer Res Treat Off J Korean Cancer Assoc* [Internet]. 2023 Apr [cited 2024 Apr 5];55(2):351–66. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10101787/>
15. Olink® Explore 3072 high-throughput proteomics platform now available: Significantly expands Olink’s protein library for biomarker discovery | Olink Holding AB [Internet]. [cited 2024 Apr 5]. Available from: <https://investors.olink.com/news-releases/news-release-details/olinkr-explore-3072-high-throughput-proteomics-platform-now/>

16. Prostate Specific Antigen (PSA) ELISA for serum or plasma 2-25ng/ml Dialab [Internet]. [cited 2024 May 2]. Available from: <https://www.alphalabs.co.uk/z00338>
17. Mosley JD, Feng Q, Wells QS, Van Driest SL, Shaffer CM, Edwards TL, et al. A study paradigm integrating prospective epidemiologic cohorts and electronic health records to identify disease biomarkers. *Nat Commun* [Internet]. 2018 Aug 30 [cited 2023 Dec 7];9:3522. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6117367/>
18. Holmes MV, Davey Smith G. Can Mendelian Randomization Shift into Reverse Gear? *Clin Chem*. 2019 Mar;65(3):363–6.
19. Balmain A. Peto’s paradox revisited: black box vs mechanistic approaches to understanding the roles of mutations and promoting factors in cancer. *Eur J Epidemiol* [Internet]. 2023 Dec 1 [cited 2024 Apr 5];38(12):1251–8. Available from: <https://doi.org/10.1007/s10654-022-00933-x>
20. Choi SW, Mak TSH, O’Reilly PF. Tutorial: a guide to performing polygenic risk score analyses. *Nat Protoc* [Internet]. 2020 Sep [cited 2023 Oct 5];15(9):2759–72. Available from: <https://www.nature.com/articles/s41596-020-0353-1>
21. Davey Smith G, Hofman A, Brennan P. Chance, ignorance, and the paradoxes of cancer: Richard Peto on developing preventative strategies under uncertainty. *Eur J Epidemiol* [Internet]. 2023 Dec 1 [cited 2024 Mar 26];38(12):1227–37. Available from: <https://doi.org/10.1007/s10654-023-01090-5>
22. Mitchell RE, Hemani G, Dudding T, Corbin L, Harrison S, Paternoster L. UK Biobank Genetic Data: MRC-IEU Quality Control, version 2, 18/01/2019.
23. Fernandez-Rozadilla C, Timofeeva M, Chen Z, Law P, Thomas M, Schmit S, et al. Deciphering colorectal cancer genetics through multi-omic analysis of 100,204 cases and 154,587 controls of European and east Asian ancestries. *Nat Genet*. 2023 Jan;55(1):89–99.
24. McKay JD, Hung RJ, Han Y, Zong X, Carreras-Torres R, Christiani DC, et al. Large-scale association analysis identifies new lung cancer susceptibility loci and heterogeneity in genetic susceptibility across histological subtypes. *Nat Genet* [Internet]. 2017 Jul [cited 2023 Oct 11];49(7):1126–32. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5510465/>
25. van Erp S, Oberski DL, Mulder J. Shrinkage priors for Bayesian penalized regression. *J Math Psychol* [Internet]. 2019 Apr 1 [cited 2024 Jan 4];89:31–50. Available from: <https://www.sciencedirect.com/science/article/pii/S0022249618300567>
26. Ge T. PRS-CS [Internet]. 2018 [cited 2023 Nov 8]. Available from: <https://github.com/getian107/PRScs>
27. Ge T, Irvin MR, Patki A, Srinivasasainagendra V, Lin YF, Tiwari HK, et al. Development and validation of a trans-ancestry polygenic risk score for type 2 diabetes in diverse populations. *Genome Med* [Internet]. 2022 Jun 29 [cited 2024 Feb 14];14(1):70. Available from: <https://doi.org/10.1186/s13073-022-01074-2>
28. Martin AR, Kanai M, Kamatani Y, Okada Y, Neale BM, Daly MJ. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat Genet* [Internet]. 2019 Apr [cited 2024 Feb 14];51(4):584–91. Available from: <https://www.nature.com/articles/s41588-019-0379-x>
29. Sun BB, Chiou J, Traylor M, Benner C, Hsu YH, Richardson TG, et al. Plasma proteomic associations with genetics and health in the UK Biobank. *Nature* [Internet]. 2023 Oct [cited 2024 Jan 5];622(7982):329–38. Available from: <https://www.nature.com/articles/s41586-023-06592-6>

30. Enroth S, Hallmans G, Grankvist K, Gyllensten U. Effects of Long-Term Storage Time and Original Sampling Month on Biobank Plasma Protein Concentrations. *eBioMedicine* [Internet]. 2016 Oct 1 [cited 2024 Feb 14];12:309–14. Available from: [https://www.thelancet.com/journals/ebiom/article/PIIS2352-3964\(16\)30394-2/fulltext](https://www.thelancet.com/journals/ebiom/article/PIIS2352-3964(16)30394-2/fulltext)
31. McCaw Z. RNOmni: Rank Normal Transformation Omnibus Test [Internet]. 2023 [cited 2024 May 22]. Available from: <https://cran.r-project.org/web/packages/RNOmni/index.html>
32. MRCIEU/metaboprep: a pipeline of metabolomics data processing and quality control [Internet]. [cited 2024 Feb 9]. Available from: <https://github.com/MRCIEU/metaboprep>
33. Haycock PC, Burgess S, Wade KH, Bowden J, Relton C, Davey Smith G. Best (but oft-forgotten) practices: the design, analysis, and interpretation of Mendelian randomization studies. *Am J Clin Nutr*. 2016 Apr;103(4):965–78.
34. Davies NM, Holmes MV, Smith GD. Reading Mendelian randomisation studies: a guide, glossary, and checklist for clinicians. *BMJ* [Internet]. 2018 Jul 12 [cited 2023 Feb 4];362:k601. Available from: <https://www.bmj.com/content/362/bmj.k601>
35. Burgess S, Thompson SG. Use of allele scores as instrumental variables for Mendelian randomization. *Int J Epidemiol* [Internet]. 2013 Aug [cited 2023 Oct 24];42(4):1134–44. Available from: <https://academic.oup.com/ije/article-lookup/doi/10.1093/ije/dyt093>
36. Davies NM, Holmes MV, Davey Smith G. Reading Mendelian randomisation studies: a guide, glossary, and checklist for clinicians. *BMJ* [Internet]. 2018 Jul 12 [cited 2023 Oct 11];k601. Available from: <https://www.bmj.com/lookup/doi/10.1136/bmj.k601>
37. Zheng J, Baird D, Borges MC, Bowden J, Hemani G, Haycock P, et al. Recent Developments in Mendelian Randomization Studies. *Curr Epidemiol Rep* [Internet]. 2017 [cited 2024 Mar 4];4(4):330–45. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5711966/>
38. Sanderson E, Glymour MM, Holmes MV, Kang H, Morrison J, Munafò MR, et al. Mendelian randomization. *Nat Rev Methods Primer* [Internet]. 2022 Feb 10 [cited 2023 Nov 10];2:6. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7614635/>
39. Sun BB, Chiou J, Traylor M, Benner C, Hsu YH, Richardson TG, et al. Genetic regulation of the human plasma proteome in 54,306 UK Biobank participants [Internet]. *bioRxiv*; 2022 [cited 2023 Oct 5]. p. 2022.06.17.496443. Available from: <https://www.biorxiv.org/content/10.1101/2022.06.17.496443v1>
40. Fauman EB, Hyde C. An optimal variant to gene distance window derived from an empirical definition of cis and trans protein QTLs. *BMC Bioinformatics* [Internet]. 2022 May 8 [cited 2023 Nov 9];23(1):169. Available from: <https://doi.org/10.1186/s12859-022-04706-x>
41. Burgess S, Thompson SG, CRP CHD Genetics Collaboration. Avoiding bias from weak instruments in Mendelian randomization studies. *Int J Epidemiol* [Internet]. 2011 Jun 1 [cited 2024 Mar 4];40(3):755–64. Available from: <https://doi.org/10.1093/ije/dyr036>
42. Burgess S, Small DS, Thompson SG. A review of instrumental variable estimators for Mendelian randomization. *Stat Methods Med Res* [Internet]. 2017 Oct [cited 2024 Mar 5];26(5):2333–55. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5642006/>

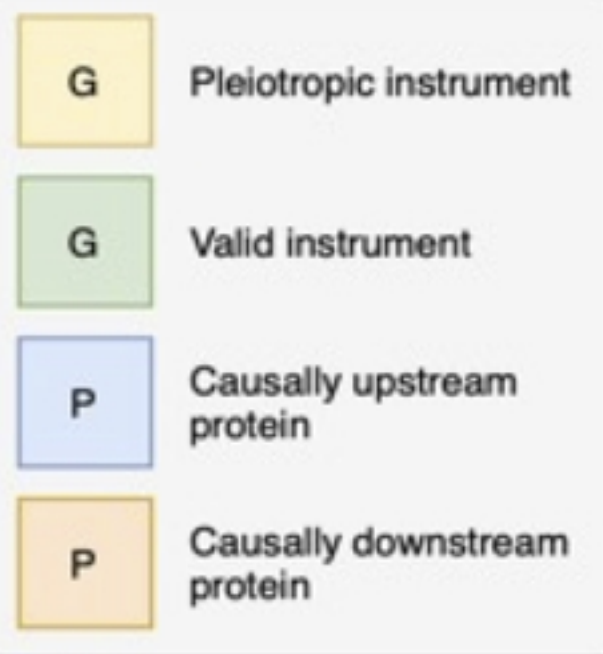
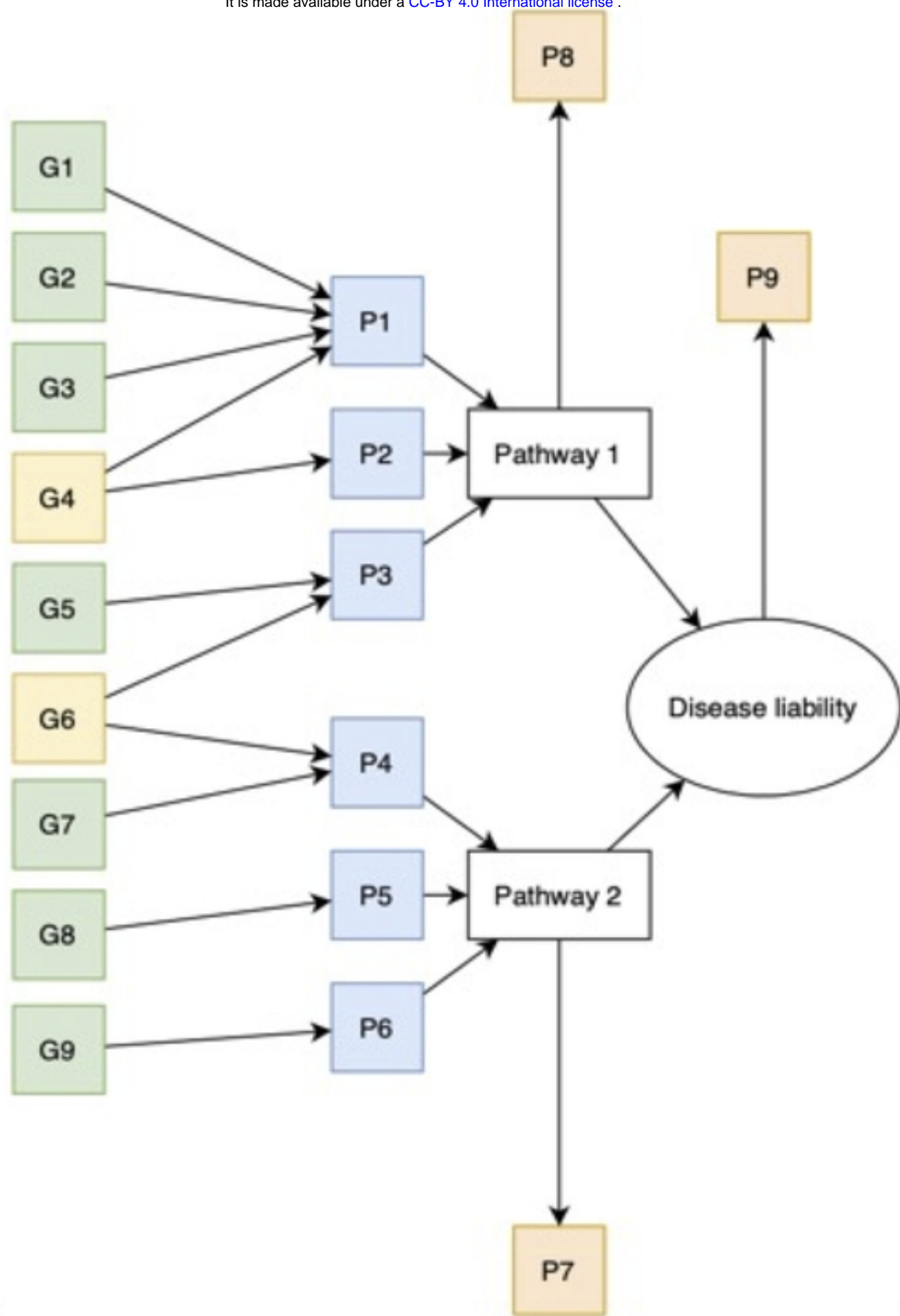
43. Burgess S, Butterworth A, Thompson SG. Mendelian Randomization Analysis With Multiple Genetic Variants Using Summarized Data. *Genet Epidemiol* [Internet]. 2013 [cited 2023 Feb 4];37(7):658–65. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1002/gepi.21758>
44. Bowden J, Hemani G, Davey Smith G. Invited Commentary: Detecting Individual and Global Horizontal Pleiotropy in Mendelian Randomization—A Job for the Humble Heterogeneity Statistic? *Am J Epidemiol* [Internet]. 2018 Dec 1 [cited 2024 May 15];187(12):2681–5. Available from: <https://doi.org/10.1093/aje/kwy185>
45. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent Estimation in Mendelian Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genet Epidemiol* [Internet]. 2016 May [cited 2024 Mar 7];40(4):304–14. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4849733/>
46. Hartwig FP, Davey Smith G, Bowden J. Robust inference in summary data Mendelian randomization via the zero modal pleiotropy assumption. *Int J Epidemiol* [Internet]. 2017 Dec 1 [cited 2024 Jan 25];46(6):1985–98. Available from: <https://doi.org/10.1093/ije/dyx102>
47. Verbanck M, Chen CY, Neale B, Do R. Detection of widespread horizontal pleiotropy in causal relationships inferred from Mendelian randomization between complex traits and diseases. *Nat Genet* [Internet]. 2018 May [cited 2023 Nov 17];50(5):693–8. Available from: <https://www.nature.com/articles/s41588-018-0099-7>
48. Zhao Q, Wang J, Hemani G, Bowden J, Small DS. Statistical inference in two-sample summary-data Mendelian randomization using robust adjusted profile score. *Ann Stat* [Internet]. 2020 Jun [cited 2024 Mar 7];48(3):1742–69. Available from: <https://projecteuclid.org/journals/annals-of-statistics/volume-48/issue-3/Statistical-inference-in-two-sample-summary-data-Mendelian-randomization-using/10.1214/19-AOS1866.full>
49. Morrison J, Knoblach N, Marcus JH, Stephens M, He X. Mendelian randomization accounting for correlated and uncorrelated pleiotropic effects using genome-wide summary statistics. *Nat Genet* [Internet]. 2020 Jul [cited 2024 Apr 19];52(7):740–7. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7343608/>
50. Foley CN, Mason AM, Kirk PDW, Burgess S. MR-Clust: clustering of genetic variants in Mendelian randomization with similar causal estimates. *Bioinforma Oxf Engl*. 2021 May 1;37(4):531–41.
51. Foley CN. [cnfoley/mrclust](https://github.com/cnfoley/mrclust) [Internet]. 2024 [cited 2024 Mar 18]. Available from: <https://github.com/cnfoley/mrclust>
52. Grant AJ, Gill D, Kirk PDW, Burgess S. Noise-augmented directional clustering of genetic association data identifies distinct mechanisms underlying obesity. *PLoS Genet*. 2022 Jan;18(1):e1009975.
53. Grant A. [aj-grant/navmix](https://github.com/aj-grant/navmix) [Internet]. 2024 [cited 2024 Mar 18]. Available from: <https://github.com/aj-grant/navmix>
54. Burgess S, Foley CN, Allara E, Staley JR, Howson JMM. A robust and efficient method for Mendelian randomization with hundreds of genetic variants. *Nat Commun* [Internet]. 2020 Jan 17 [cited 2024 Mar 18];11(1):376. Available from: <https://www.nature.com/articles/s41467-019-14156-4>
55. Darrous L, Hemani G, Davey Smith G, Kutalik Z. PheWAS-based clustering of Mendelian Randomisation instruments reveals distinct mechanism-specific causal effects between obesity and educational attainment. *Nat Commun*. 2024 Feb 15;15(1):1420.

56. Hu X, Cai M, Xiao J, Wan X, Wang Z, Zhao H, et al. Benchmarking Mendelian Randomization methods for causal inference using genome-wide association study summary statistics [Internet]. medRxiv; 2024 [cited 2024 Apr 18]. p. 2024.01.03.24300765. Available from: <https://www.medrxiv.org/content/10.1101/2024.01.03.24300765v1>
57. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, et al. The MR-Base platform supports systematic causal inference across the human phenome. Loos R, editor. eLife [Internet]. 2018 May 30 [cited 2023 Feb 4];7:e34408. Available from: <https://doi.org/10.7554/eLife.34408>
58. Foley CN, Staley JR, Breen PG, Sun BB, Kirk PDW, Burgess S, et al. A fast and efficient colocalization algorithm for identifying shared genetic risk factors across multiple traits. Nat Commun [Internet]. 2021 Feb 3 [cited 2024 Feb 14];12(1):764. Available from: <https://www.nature.com/articles/s41467-020-20885-8>
59. Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, et al. Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using Summary Statistics. PLoS Genet [Internet]. 2014 May 15 [cited 2023 Nov 1];10(5):e1004383. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4022491/>
- 60.: Data-Field 40005 [Internet]. [cited 2024 Feb 15]. Available from: <https://biobank.ndph.ox.ac.uk/showcase/field.cgi?id=40005>
- 61.: Data-Field 3166 [Internet]. [cited 2024 Feb 15]. Available from: <https://biobank.ndph.ox.ac.uk/showcase/field.cgi?id=3166>
62. deCODE genetics [Internet]. 2012 [cited 2023 Dec 1]. SCIENCE. Available from: <https://www.decode.com/research/>
63. Riboli E, Hunt KJ, Slimani N, Ferrari P, Norat T, Fahey M, et al. European Prospective Investigation into Cancer and Nutrition (EPIC): study populations and data collection. Public Health Nutr. 2002 Dec;5(6B):1113–24.
64. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, et al. The MR-Base platform supports systematic causal inference across the human phenome. eLife. 2018 May 30;7:e34408.
65. Ethics [Internet]. [cited 2024 Jan 4]. Available from: <https://www.ukbiobank.ac.uk/learn-more-about-uk-biobank/about-us/ethics>
- 66.: External Info : returning_results [Internet]. [cited 2024 Jan 5]. Available from: https://biobank.ndph.ox.ac.uk/showcase/exinfo.cgi?src=returning_results

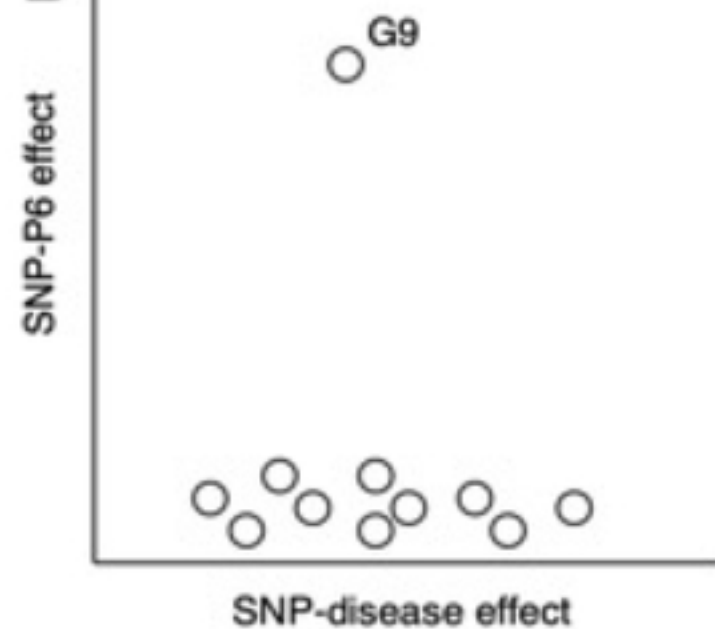
A

General population

medRxiv preprint doi: <https://doi.org/10.1101/2024.10.18.24315725>; this version posted October 18, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



B



C

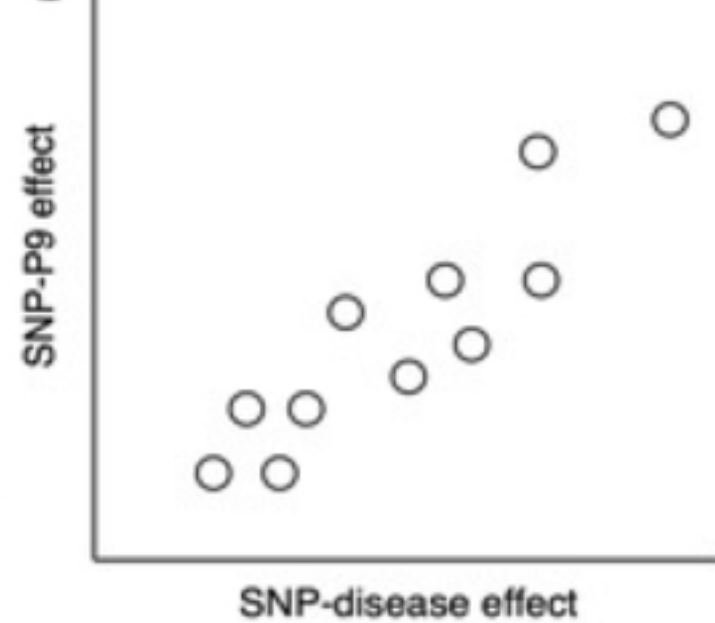
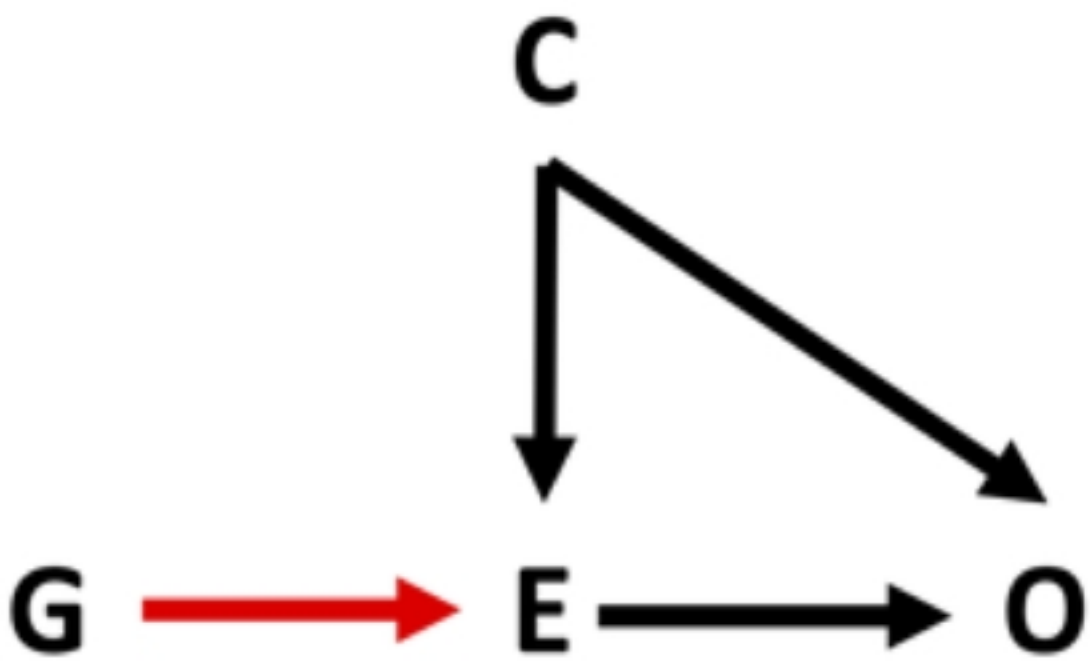
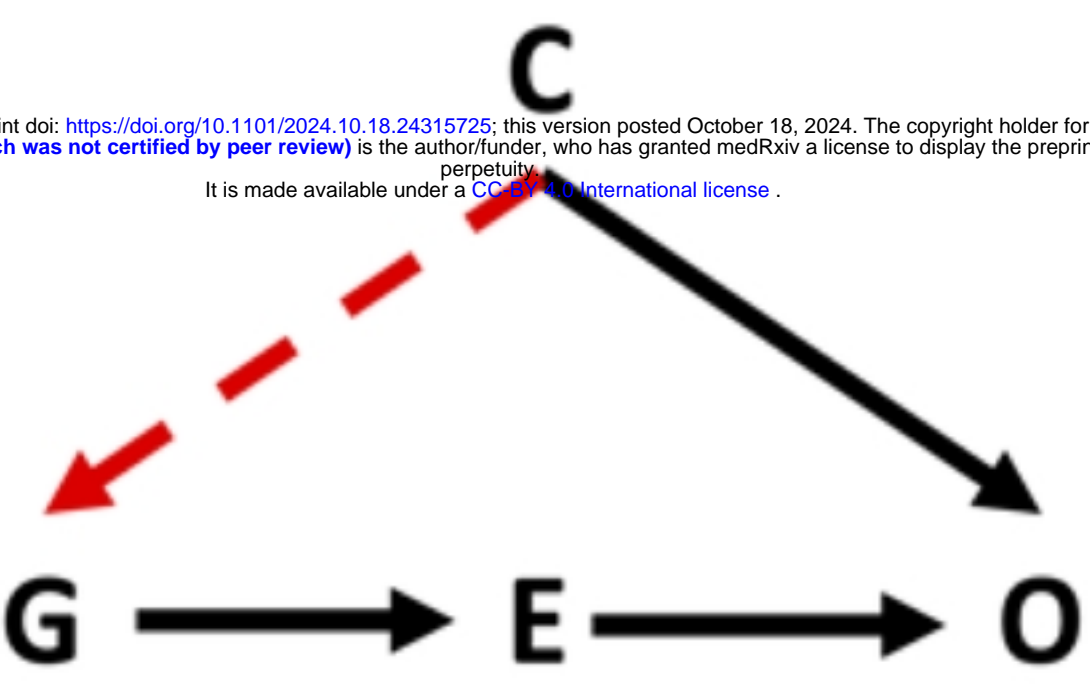


Figure 1

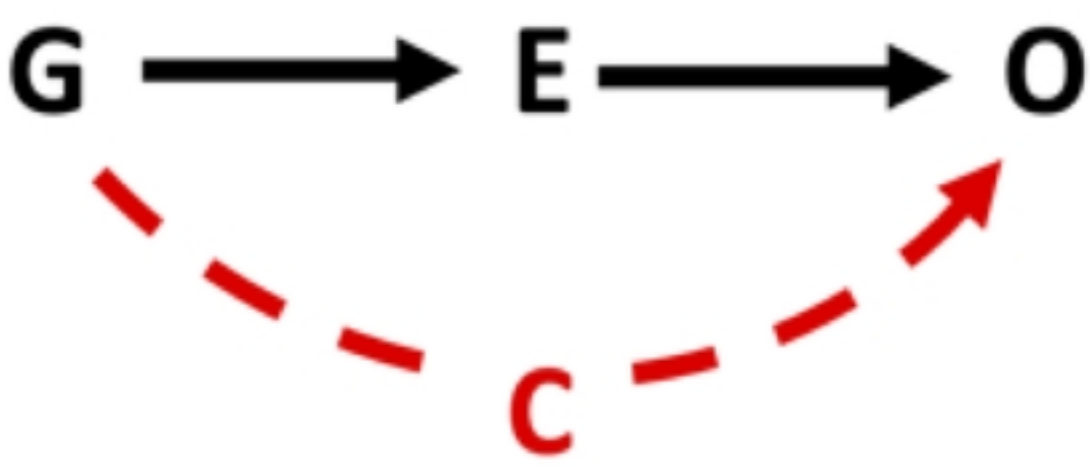


IV2

medRxiv preprint doi: <https://doi.org/10.1101/2024.10.18.24315725>; this version posted October 18, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



IV3



G – Instrumental variable (proxy)
E – Exposure
O – Outcome
C – Confounder

Figure 2