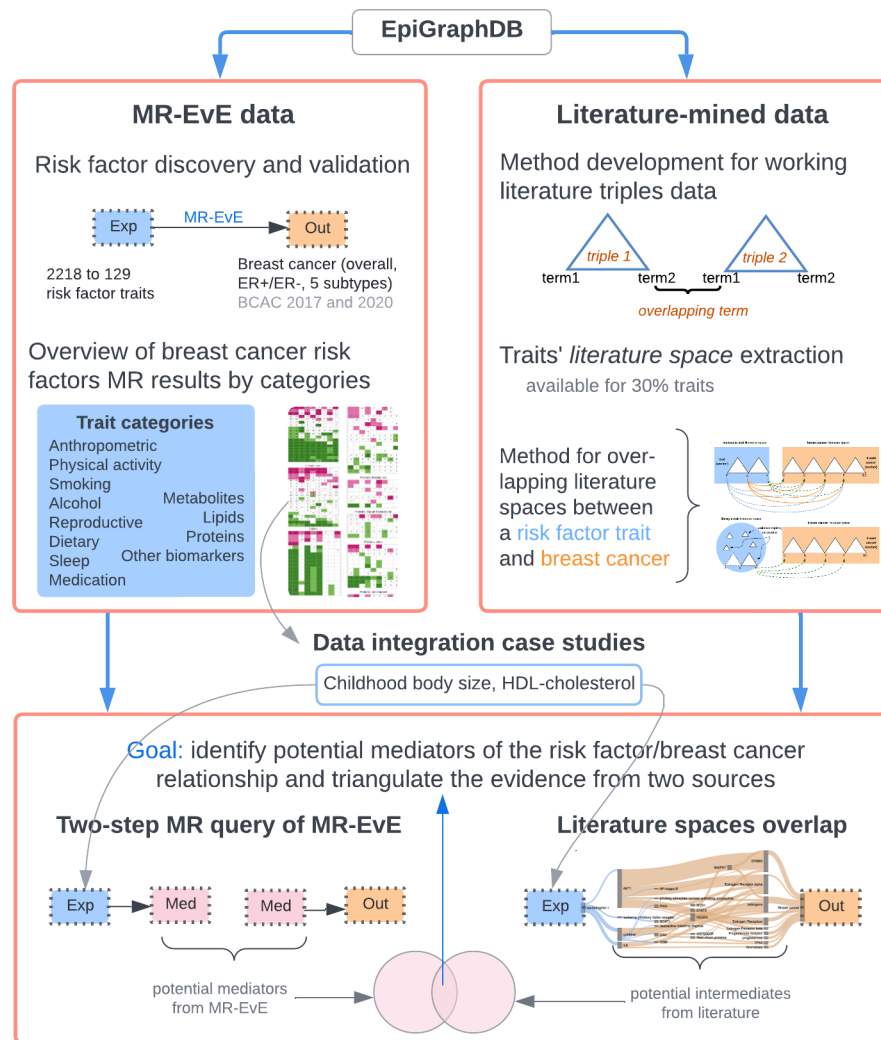


## 1 Graphical Abstract

### 2 Integrating Mendelian randomization and literature-mined evidence for breast cancer risk factors

3  
4 Marina Vabistsevits, Tim Robinson, Ben Elsworth, Yi Liu, Tom R Gaunt



5 Integrating Mendelian randomization and literature-mined  
6 evidence for breast cancer risk factors

7 Marina Vabistsevits<sup>a,b,d</sup>, Tim Robinson<sup>a,b</sup>, Ben Elsworth<sup>a,b,c</sup>, Yi Liu<sup>a,b</sup>, Tom R  
8 Gaunt<sup>a,b</sup>

<sup>a</sup>Medical Research Council Integrative Epidemiology Unit at the University of Bristol, University of  
Bristol, Oakfield 7 House, Oakfield Grove, Bristol, BS8 2BN, UK

<sup>b</sup>Population Health Sciences, Bristol Medical School, University of Bristol, Oakfield 7 House,  
Oakfield Grove, Bristol, BS8 2BN, UK

<sup>c</sup>Our Future Health, 2 New Bailey, 6 Stanley Street, Manchester, M3 5GS, UK

<sup>d</sup>University of Exeter Medical School, University of Exeter St Luke's Campus, 79 Heavitree  
Rd, Exeter, EX2 4TH, UK

<sup>e</sup>corresponding author: Marina Vabistsevits ([marina.vabistsevits@bristol.ac.uk](mailto:marina.vabistsevits@bristol.ac.uk))

---

9 **Abstract**

10 *Objective:* An increasing challenge in population health research is efficiently  
11 utilising the wealth of data available from multiple sources to investigate disease  
12 mechanisms and identify potential intervention targets. The use of biomedical  
13 data integration platforms can facilitate evidence triangulation from these differ-  
14 ent sources, improving confidence in causal relationships of interest. In this work,  
15 we aimed to integrate Mendelian randomization (MR) and literature-mined evi-  
16 dence from the EpiGraphDB biomedical knowledge graph to build a comprehen-  
17 sive overview of risk factors for developing breast cancer.

18 *Methods:* We utilised MR-EvE (“Everything-vs-Everything”) data to identify  
19 candidate risk factors for breast cancer and generate hypotheses for potential me-  
20 diators of their effect. We also integrated this data with literature-mined relation-  
21 ships, which were extracted by overlapping literature spaces of risk factors and  
22 breast cancer. The literature-based discovery (LBD) results were followed up by  
23 validation with two-step MR to triangulate the findings from two data sources.

24 *Results:* We identified 129 novel and established lifestyle risk factors and  
25 molecular traits with evidence of an effect on breast cancer, and made the MR  
26 results available in an R/Shiny app ([https://mvab.shinyapps.io/MR\\_](https://mvab.shinyapps.io/MR_heatmaps/)  
27 [heatmaps/](https://mvab.shinyapps.io/MR_heatmaps/)). We developed an LBD approach for identifying potential mech-  
28 anistic intermediates of identified risk factors. We present the results of MR and  
29 literature evidence integration for two case studies (childhood body size and HDL-  
30 cholesterol), demonstrating their complementary functionalities.

31 *Conclusion:* We demonstrate that MR-EvE data offers an efficient hypothesis-

32 generating approach for identifying disease risk factors. Moreover, we show that  
33 integrating MR evidence with literature-mined data may be used to identify causal  
34 intermediates and uncover the mechanisms behind the disease.

35 *Keywords:* mendelian randomization, genome-wide association study, causal  
36 inference, breast cancer, literature-based discovery, knowledge graph, literature  
37 mining

---

## 38 1. Introduction

39 Triangulation, a process of integrating evidence from multiple methodologies,  
40 has become increasingly important in population health research [1, 2]. Differ-  
41 ent evidence sources are likely to have unique and unrelated sources of bias, so  
42 consistent evidence from such sources improves confidence in the causal relation-  
43 ship between a risk factor and outcomes of interest [3]. A potential obstacle to the  
44 widespread adaptation of triangulation in epidemiological research is the challenge  
45 of integrating and harmonising the large volume of data from different datasets.  
46 Multiple platforms [4, 5, 6, 7] have been developed to facilitate access to com-  
47 bined and curated biomedical datasets. Using these data integration platforms will  
48 facilitate the study of disease aetiology.

49 EpiGraphDB (<https://epigraphdb.org>) [4] is a biomedical knowledge  
50 graph that was designed to facilitate data mining of epidemiological relationships.  
51 The unique advantages of this platform are the availability of both Mendelian ran-  
52 domization (MR) [8, 9] and literature-mined relationships between molecular and  
53 lifestyle traits and disease outcomes, providing a valuable combination of data to  
54 systematically investigate causal relationships.

55 MR is an approach to causal inference that uses genetic variants as instru-  
56 mentable variables (IVs) to infer whether a modifiable health exposure influences  
57 a disease outcome [8]. Studying disease causality with MR methods has been  
58 growing in popularity over the past decade and a half [9, 10], with the increasing  
59 availability of data from genome-wide association studies (GWAS) [11, 12] and  
60 well-developed analysis frameworks [13]. EpiGraphDB provides access to MR-  
61 EvE (Mendelian randomization “Everything-vs-Everything”) pairwise causal esti-  
62 mates [14] between thousands of traits from the OpenGWAS database (<https://gwas.mrcieu.ac.uk>) [11]. MR-EvE data can be used to explore poten-  
63 tial causal links between thousands of traits for knowledge discovery or hypothesis  
64 generation.  
65

66 Another aspect of knowledge discovery is extracting data from the biomed-  
67 ical literature, also known as literature-based discovery (LBD) [15]. In addition to  
68 MR-EvE, EpiGraphDB also contains literature-mined relationships extracted from

69 publications in PubMed [16] and mapped to specific phenotypes in the MR-EvE  
70 data. LBD entails connecting information that is explicitly stated in multiple pub-  
71 lications across different research disciplines aiming to deduce connections that  
72 have not been explicitly stated [17, 18], offering the potential to identify mechanis-  
73 tic pathways or mediators for causal relationships identified using MR-EvE.

74 In this work, we use EpiGraphDB to build a comprehensive picture of the aeti-  
75 ology of breast cancer by combining evidence from MR and literature-mined data.  
76 Breast cancer is a useful exemplar, as it is a heterogeneous disease with a complex  
77 aetiology, affected by both genetic and lifestyle factors [19, 20, 21, 22] that also  
78 has a large body of published research evidence, suitable for data mining. We seek  
79 to generate new mechanistic hypotheses for causal relationships using evidence in-  
80 tegration to inform prevention and intervention strategies that could reduce disease  
81 burden.

82  
83

### Statement of Significance

<b>Problem or Issue</b>	Limited integration of multiple sources of evidence to enhance causal inference and identify mechanisms in epidemiological research
<b>What is already known</b>	While MR studies are abundant, few utilise complementary biomedical data sources to enhance causal evidence or generate new hypotheses for underlying mechanisms
<b>What this paper adds</b>	By utilising MR-EvE and literature-mined data in EpiGraphDB, this study introduces a systematic, integrated approach for causal inference, applied to breast cancer risk factors. It demonstrates how combining MR and LBD enhances the identification of mediating traits, providing a model that can be applied to other complex diseases

## 84 2. Methods

### 85 2.1. Study Design

86 We present a summary of our EpiGraphDB data-focused study in Figure 1.  
87 First, we explored MR-EvE causal estimates between various exposure traits and  
88 breast cancer to generate a list of potential causal risk factors and biomarkers.  
89 Then, we used literature-mined relationships from EpiGraphDB to dissect the risk  
90 factor/breast cancer relationships. We designed an LBD method for overlapping lit-  
91 erature spaces of traits to identify potential intermediates between individual traits  
92 and breast cancer.

93 We demonstrate evidence integration between MR and LBD on two case study  
94 breast cancer risk factors: *childhood body size* and *HDL-cholesterol*. We queried

95 MR-EvE data using a two-step MR approach [23, 24] to identify potential medi-  
96 ators (i.e. traits affected by each risk factor that in turn affect breast cancer risk).  
97 We also determined potential intermediates using LBD and evaluated them using  
98 MR and data from OpenGWAS. In this way, we triangulated evidence for the iden-  
99 tified mediators from two independent data sources. Analysis code is available at  
100 <https://github.com/mvab/epigraphdb-breast-cancer> and [https://github.com/mvab/epigraphdb\\_mr\\_literature\\_queries](https://github.com/mvab/epigraphdb_mr_literature_queries).  
101

## 102 2.2. *EpiGraphDB*

103 EpiGraphDB (<https://epigraphdb.org>) is a graph database that inte-  
104 grates biomedical and epidemiological data, built to support data mining of risk  
105 factor/disease relationships. It contains trait relationships (observational and ge-  
106 netic), literature-mined relationships, biological pathways, protein-protein interac-  
107 tions (PPIs), drug–target relationships and other data sources [4]. EpiGraphDB is  
108 integrated with the OpenGWAS database data [11] (<https://gwas.mrcieu.ac.uk>),  
109 providing access to GWAS studies via the GWAS node of the graph.  
110 The number of GWAS studies in OpenGWAS at the time of writing is 50,040; of  
111 these, 34,494 are integrated within EpiGraphDB. EpiGraphDB (version 0.3.0) was  
112 accessed via R package *epigraphdb-r* (version 0.2.3).

## 113 2.3. *Breast cancer data*

114 EpiGraphDB provides access to multiple breast cancer GWAS datasets de-  
115 posited in OpenGWAS. We used the most recent releases of the Breast Cancer As-  
116 sociation Consortium (BCAC) (<http://bcac.ccge.medschl.cam.ac.uk/>)  
117 meta-analyses conducted in 2017 (BCAC 2017, N=228,951) [25] and 2020 (BCAC  
118 2020, N=247,173) [26]. Since only BCAC 2017 is available in EpiGraphDB/OpenGWAS,  
119 this dataset was used for the MR-EvE risk factor discovery analysis. Both BCAC  
120 2017 and 2020 were utilized for the validation stage. In addition to the overall sam-  
121 ple, the datasets contain the ER+ and ER- sub-samples and five molecular subtypes  
122 (Supplementary Table 1). The study groups in the BCAC cohort do not include UK  
123 Biobank or other trait cohorts used in this study. The cohort includes European-  
124 ancestry individuals only.

## 125 2.4. *Mendelian randomization and MR-EvE data*

126 MR is an application of instrumental variable analysis where genetic variants  
127 are used as proxies (genetic instruments) to estimate the causal relationship be-  
128 tween a modifiable health exposure and a disease outcome [8, 9]. MR-EvE (MR  
129 “Everything-vs-Everything”) data is stored as bidirectional relationships between  
130 GWAS traits in EpiGraphDB, available for 11,789 GWAS traits at the time of writ-  
131 ing. The MR-EvE estimates were generated using the MR Mixture-of-Experts

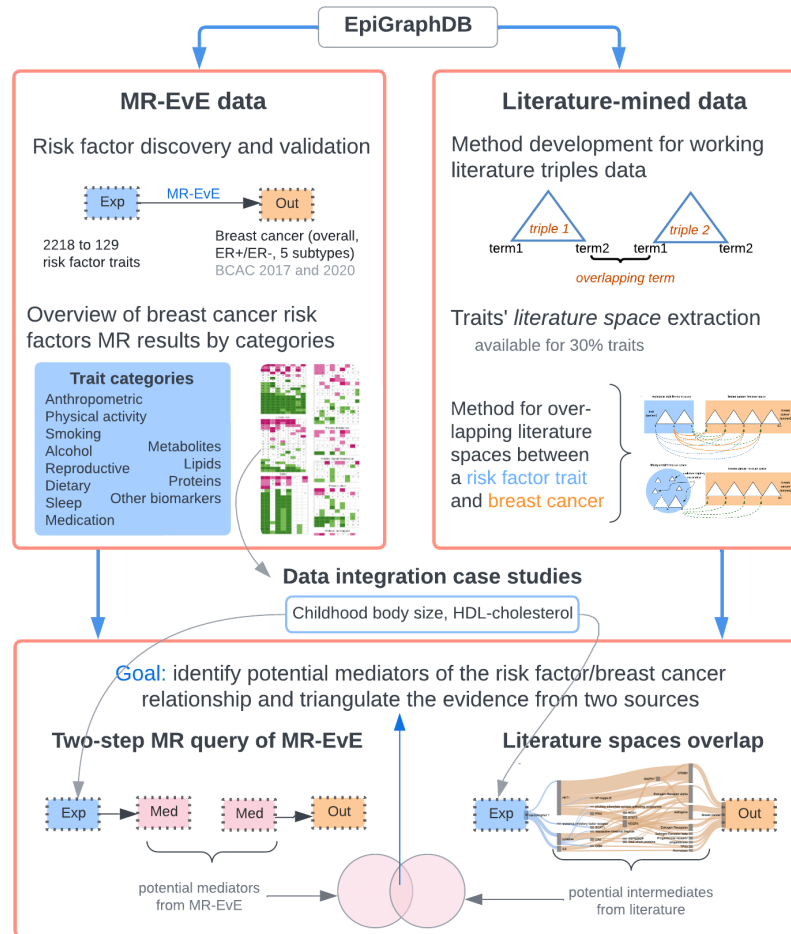


Figure 1: **Study overview.** Employing EpiGraphDB, MR-EvE data was used to identify breast cancer risk factors, and literature-mined data was used to design a literature-based discovery (LBD) method for identifying potential intermediates between risk factors and breast cancer. For two case studies, we present evidence integration between MR-EvE-guided and LBD-identified potential mediators. Abbreviations/definitions: MR-EvE – Mendelian Randomization “Everything-vs-Everything”; Exp – exposure; Med – mediator; Out – outcome; BCAC – Breast Cancer Association Consortium (refers to meta-analyses breast cancer cohort GWAS data); literature space – a collection of literature-mined relationships between pairs of biological terms related to a phenotype.

132 (MR-MoE) method by Hemani *et al* [14] before integration into EpiGraphDB. MR-  
 133 MoE is a machine learning framework that automated the selection of instrument  
 134 SNPs and the MR method for use in any specific causal analysis by predicting  
 135 the model of pleiotropy and relating that prediction to the most appropriate model

136 [14]. The method selects the MR ‘best estimate’ of causal relationships between  
137 the exposures and the outcome traits.

138 The interconnection of many traits in the MR-EvE data makes it a valuable  
139 resource for identifying potential mediators and confounders [27] for a causal rela-  
140 tionship of interest (<https://epigraphdb.org/confounder>) [4]. In this  
141 work, we make use of this functionality to generate hypotheses for breast cancer  
142 risk factors and potential mediators.

## 143 2.5. Identifying breast cancer risk factor traits from MR-EvE

### 144 2.5.1. Discovery

145 We queried EpiGraphDB to extract MR-EvE results for all available GWAS  
146 traits (N = 2218, Supplementary Table 2) with breast cancer as the outcome. Figure  
147 2 summarises the discovery and validation stages of the analysis. The initial set of  
148 traits linked to breast cancer with any MR estimate was reduced to traits that passed  
149 False Discovery Rate (FDR), taking traits with the corrected p-value < 0.05 (N =  
150 378) into the downstream analysis (Figure 2).

151 The traits were grouped into 12 categories, which included widely-known can-  
152 cer risk factor categories (anthropometric traits, reproductive traits, physical ac-  
153 tivity, smoking, alcohol), categories with a growing research interest (dietary in-  
154 take traits, sleep traits, medication), and molecular traits (metabolites, lipids, pro-  
155 teins, and other biomarkers) (N = 202). The traits outside of these categories were  
156 grouped as ‘Other’ and excluded. The full data from the discovery stage is avail-  
157 able in Supplementary Table 3.

### 158 2.5.2. Validation

159 The results of the MR-EvE rapid screen of risk factors were validated by re-  
160 analysing the identified relationships using the robust practices in two-sample MR  
161 [28], because MR-EvE ‘best estimates’ selected by the MR-MoE method do not  
162 always correspond to the standard methods used in MR studies [14]. The validation  
163 results in this study are reported as per the STROBE-MR guidelines [29, 30].

164 For the majority of exposure traits, instruments were extracted using the genome-  
165 wide approach, i.e. SNPs with p-value <  $5 \times 10^{-8}$  and clumped with  $r^2 < 0.001$   
166 using 1000 Genomes European LD reference panel. For protein traits, *cis*-region  
167 instruments were prioritised (within 1Mb on either side of the protein-encoding  
168 gene) [31]. The extracted *cis*-SNPs were filtered to p-value <  $5 \times 10^{-8}$  or a less  
169 stringent ‘local’ threshold of p-value < 0.05/nSNPs, where nSNPs is the number  
170 of SNPs in the *cis* region. The clumping was done using  $r^2 < 0.01$ , or  $r^2 < 0.001$   
171 if the higher threshold produced  $\geq 7$  SNPs. If no instruments were available under  
172 these criteria, genome-wide extraction of instruments was applied. The details of

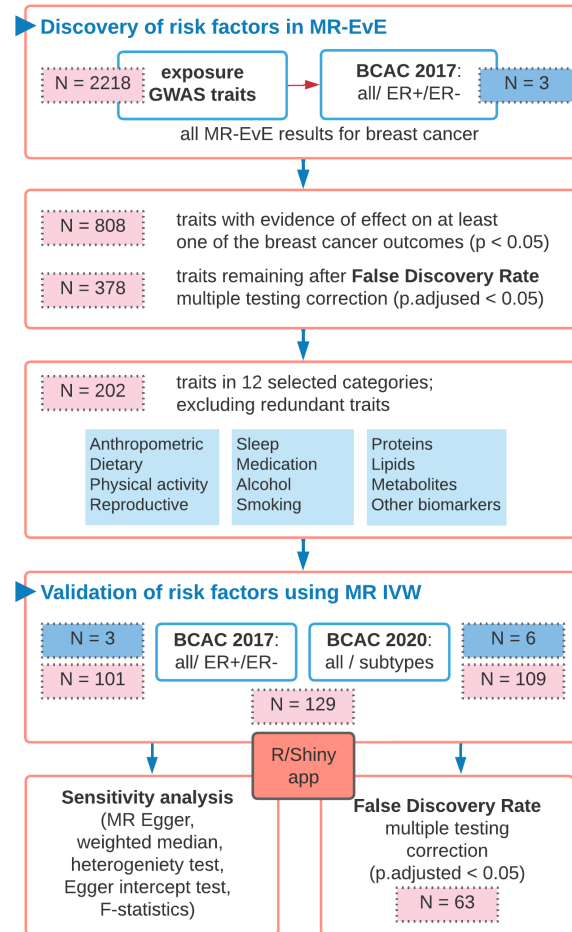


Figure 2: **MR-EvE evidence processing workflow, including Discovery and Validation stages.** The boxes show the number of exposures (pink) and breast cancer outcomes (blue) after/in each filtering step; MR-EvE results validation Shiny app: [https://mvab.shinyapps.io/MR\\_heatmaps/](https://mvab.shinyapps.io/MR_heatmaps/). BCAC - Breast Cancer Association Consortium

173 instrument extraction for protein traits are provided in Supplementary Tables 4 and  
174 5.

175 The validation analysis was performed on breast cancer molecular subtypes  
176 data (BCAC 2020) in addition to the overall sample and ER+/ER- subtypes (BCAC  
177 2017). The inverse-variance weighted (IVW) [32] approach was used as the main  
178 two-sample MR method. To correct for multiple testing in the validation step, the  
179 FDR cut-off of 0.05 for the corrected p-value was applied to identify results with



180 more robust evidence of effect (Supplementary Table 4) (Figure 2).

181 To investigate potential violations of the MR assumptions and validate the ro-  
182 bustness of MR results, additional analyses were performed. This included pleiotropy-  
183 robust sensitivity analyses such as MR-Egger [33] and weighted median [34]. The  
184 Egger intercept test [33] was used to test for directional horizontal pleiotropy, and  
185 Cochran’s Q statistic [35] was calculated to quantify heterogeneity among SNPs,  
186 which is indicative of potential pleiotropy. F-statistic was used to assess instrument  
187 strength [36] (Supplementary Tables 4 and 6). Sensitivity analysis results were not  
188 used for additional filtering of all identified risk factor traits to avoid information  
189 loss due to dichotomisation [37], but were taken into account when evaluating in-  
190 dividual traits. Additionally, the sensitivity results were included in the MR results  
191 heatmap visualisation to provide context on the performance in sensitivity tests for  
192 all traits. All MR and sensitivity analyses were performed using the *TwoSampleMR*  
193 R package (version 0.5.6) [13].

## 194 2.6. Using literature-mined relationships in *EpiGraphDB*

195 *EpiGraphDB* contains literature-mined relationships between pairs of terms  
196 (also known as ‘literature triples’) derived from the published literature [38]. The  
197 underlying literature data comes from SemMedDB (Semantic MEDLINE Database)  
198 [16], a well-established repository of literature-mined semantic triples, i.e. ‘*subject*  
199 *term 1 – predicate – object term 2*’, mined from titles and abstracts of nearly 30  
200 million biomedical articles in PubMed, using SemRep [39, 40] (further details are  
201 available in Appendix A). A subject/object term can be any biological entity (gene  
202 name, drug, phenotype, disease); a predicate is a verb that represents a relationship  
203 between the two terms (*affects, causes, inhibits, reduces, associated with, etc.*).

### 204 2.6.1. Literature space

205 Literature triples of biomedical terms are mapped to specific GWAS traits in  
206 *EpiGraphDB*. In this work, we introduce the concept of a ‘*literature space*’. The  
207 literature space for a trait (e.g. breast cancer) comprises a set of semantic triples  
208 extracted from the published literature for this trait. Literature space can be consid-  
209 ered an ‘automated literature overview’ of conceptual relationships between bio-  
210 logical entities related to the search term, i.e. breast cancer. In the literature space,  
211 each relationship triple has a score that indicates how frequently it occurs in the  
212 published literature within a defined time frame (linked to PubMed IDs), demon-  
213 strating how well-established this relationship is.

### 214 2.6.2. Literature spaces overlap method

215 In the field of LBD, the underlying principle of working with triple-based liter-  
216 ature data is known as Swanson linking [41, 42]. In this approach, two separately

217 published results (e.g. A-B and B-C relationships) are combined to identify evi-  
218 dence for the A-C relationship that is unknown or unexplored. The discovery can  
219 be open and closed: in open, A is given to identify Bs linked to Cs, generating hy-  
220 potheses; in closed, A and C are fixed, and Bs that link them are identified, testing  
221 hypotheses or finding intermediates [43, 18] (details in Appendix A).

222 In this work, we present an extension to the simple triple linking method used  
223 in MELODI-Presto [38], which focused on the overlap of two triples (i.e. A-B  
224 and B-C only) (3a). Here, we combine multiple triples from two literature spaces  
225 (exposure trait and breast cancer) into a so-called '*literature overlap*'. We apply  
226 multi-step triple linking within and between the literature spaces to generate hy-  
227 potheses for the potential mediators of the exposure trait's effect on breast cancer.  
228 The literature data pre-processing steps and the method details are provided in Ap-  
229 pendix A. The overview of our LBD method is presented in Figure 3, including  
230 closed (3b) and open (3c) discovery, and an example application (3d).

## 231 2.7. MR and literature evidence integration

232 We integrated the evidence from MR-EvE and literature mining for potential  
233 mediators of two breast cancer risk factors that we selected as case studies: *child-*  
234 *hood body size* and *HDL-cholesterol*.

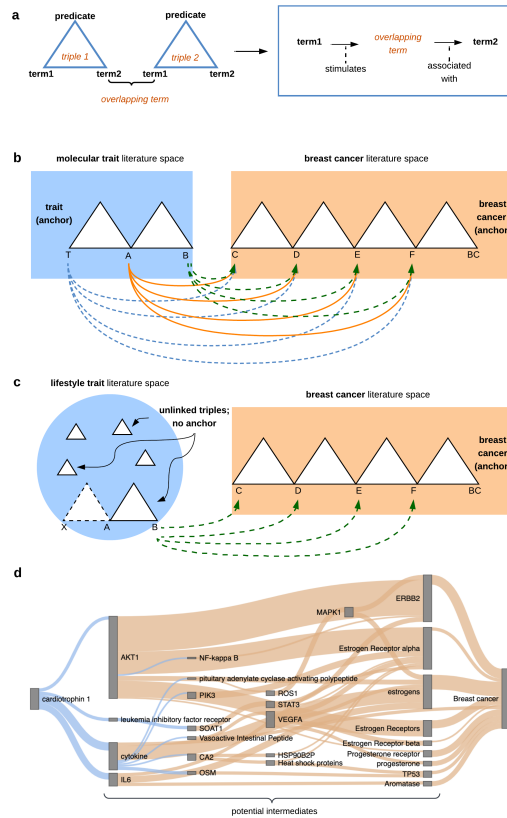
### 235 2.7.1. Potential mediators from MR-EvE data

236 We queried the MR-EvE data to identify potential mediators of each case study  
237 risk factor using the two-step MR framework [23, 24, 44]. The mediator search re-  
238 sults were corrected for multiple testing using FDR at the individual case study  
239 level, for each step within two-step MR: (*step 1*) a risk factor trait effect on poten-  
240 tial mediators, and (*step 2*) potential mediators' effects on breast cancer outcomes  
241 (BCAC 2017 data only).

242 The mediators passing the FDR-adjusted p-value  $< 0.05$  filtering in both steps  
243 and restricted to the earlier specified 12 trait categories were taken into the valida-  
244 tion step. The validation for two-step MR results included the same analysis and  
245 sensitivity methods used for the main risk factors (2.5.2). The validated two-step  
246 MR estimates were also corrected for FDR.

### 247 2.7.2. Potential mediators from the LBD literature overlap approach

248 The literature spaces for risk factor traits and breast cancer were constructed  
249 using the literature data linked to the following GWAS traits: breast cancer (com-  
250 bination of two traits: '*breast cancer*' - *ieu-a-1126* and '*malignant neoplasm of*  
251 *the breast*' - *finn-a-C3\_BREAST*), HDL ('*HDL-cholesterol*' - *ukb-d-30760\_irnt*),  
252 childhood body size ('*childhood obesity*' - *ieu-a-1096*) (Supplementary Tables 14-  
253 15).



**Figure 3: Schemas of the literature-based discovery method, demonstrating the approach behind literature spaces overlap. (a)** A basic example of linking two triples via an overlapping term. **(b) Closed discovery:** overlap of two traits' literature spaces (blue - risk factor, orange - breast cancer) by chaining triples of terms between them. The outer terms are restricted ('i.e. 'anchored') to the term representing the trait (T: trait) and the term representing breast cancer outcome (BC: breast cancer). All intermediate terms are arbitrarily named A-F. The arrows represent all possible paths from T to BC, which are alternatives to the full path going through all intermediates. Closed discovery is applied to molecular traits. **(c) Open discovery:** Applied to non-molecular/lifestyle traits, where the trait is not available as a term and therefore cannot be anchored. The unlinked triples in the non-molecular trait space are first matched (via term B) to link with any triples in the breast cancer space. Then, the linked triples are connected to other triples in the non-molecular trait space (via term A), expanding the literature spaces overlap. **(d) Example of literature spaces overlap visualisation:** Sankey diagram showing how triples within and between a risk factor trait (here, cardiotrophin-1) and breast cancer spaces interconnect. The line width of each term relationship indicates how common it is in the literature (frequency score). Cardiotrophin-1 is presented here as an example. More examples, specifically the case study traits in this study, are available to view at [https://mvab.shinyapps.io/literature\\_overlap\\_sankey/](https://mvab.shinyapps.io/literature_overlap_sankey/)

254 In the literature overlap generated from triples in the risk factor trait and breast  
255 cancer literature spaces, all intermediate biological terms were evaluated as poten-  
256 tial mediators. Each term was searched in OpenGWAS for representative GWAS  
257 data to be used in MR analysis, and genetic instruments for the identified traits  
258 were extracted as previously described (2.5.2).

259 We performed MR bidirectionally between the risk factor trait and each inter-  
260 mediate to identify the direction of effect, as it may not be correct in literature-  
261 mined relationships. Then, the intermediate traits were evaluated for an effect on  
262 breast cancer. The results were validated using the analysis and sensitivity methods  
263 described previously (2.5.2), and adjusted using FDR.

264 The evidence in both steps of the two-step MR analysis provided suggestive  
265 evidence of potential mediation between the risk factor and breast cancer via a trait  
266 identified by the LBD method.

### 267 3. Results

#### 268 3.1. Breast cancer risk factors discovered from MR-EvE

269 We identified 129 traits across 12 categories that showed evidence of an effect  
270 on breast cancer risk from MR-EvE data (Figure 2). The discovery effect estimates  
271 were validated using the IVW MR analysis with sensitivity analyses, using 9 breast  
272 cancer outcome GWAS from two releases of BCAC data. The instruments were ex-  
273 tracted as described in 2.5.2. To correct for multiple testing of many exposures and  
274 outcomes in the validation step, the FDR correction was applied. The validation  
275 results and sensitivity analyses are available in Supplementary Tables 4 and 6.

276 The MR results for the final set of traits (129 traits or 63 after FDR correc-  
277 tion) are presented in Figure 4. The figure also provides an overview of sensi-  
278 tivity analysis results: potential pleiotropy “P”, heterogeneity “H”, and “X” for  
279 exposures with a low number of instruments ( $\leq 2$ ) where most sensitivity tests  
280 were not possible. The full results (202 traits) are available in the R/Shiny app  
281 ([https://mvab.shinyapps.io/MR\\_heatmaps/](https://mvab.shinyapps.io/MR_heatmaps/)) and Supplementary Ta-  
282 ble 4. The abbreviations for traits/biomarkers used throughout the section are de-  
283 fined in Supplementary Table 7. In the sections below, the results of several trait  
284 categories are summarised.

##### 285 3.1.1. Non-molecular traits

286 The majority of anthropometric traits (e.g. *body mass index*, *obesity*, *waist-to-*  
287 *hip ratio*, *childhood body size*) (Figure 4a) have consistent inverse (risk-decreasing)  
288 effects across most breast cancer subtypes. Interestingly, *impedance of whole body*  
289 has a positive (risk-increasing) effect on most outcomes, potentially capturing a

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



Figure 4: Mendelian randomization effect direction heatmaps for (A) lifestyle and anthropometric traits, (B) metabolites and lipids, (C) proteins (N=129) as exposures, and breast cancer (BCAC 2017 and 2020 overall samples and subtypes) as the outcomes. The presented traits have evidence of an effect on at least one outcome, but the full set of traits analysed in the validation step (N=202) is available in the interactive version of the plot in the app: [https://mvab.shinyapps.io/MR\\_heatmaps/](https://mvab.shinyapps.io/MR_heatmaps/) and in Supplementary Table 4. [continued on the next page]

Figure 4: [continued] The proportions of the total size of each group displayed here are: Anthropometric – 95%, Lifestyle – 57%, Lipids – 100%, Metabolites – 68%, Proteins – 54%. The colours represent effect direction: pink – positive (causal), green – negative (protective), white – null effect (based on 95% CIs of odds of breast cancer per standard deviation increase in the exposure). The darker colour shade indicates the results that passed the FDR correction. The letters summarise arbitrarily dichotomised sensitivity analysis results: "P" – potential pleiotropy ( $abs(egger\ intercept) > 0.05$ ); "H" – heterogeneity (Q-stat p-value < 0.05); "X" – low number of instruments ( $\leq 2$ ) where sensitivity analyses were not possible; an empty cell indicates that no H/P/X issues were identified. For proteins, where possible, *cis*-only instruments were used (marked with an asterisk \*). The proteins were grouped by major pathways defined in the Reactome database [45] for display purposes. Proteins are displayed by the encoding gene name (defined in Supplementary Table 7). The exposure traits are from mixed-sex samples unless otherwise specified (F: female-only) or are female-specific reproductive traits. The details on breast cancer outcome samples are provided in Supplementary Table 1.

290 different aspect of adiposity. The traits in this group show evidence of high instru-  
291 ment heterogeneity ('H' in Figure 4a) based on the arbitrary threshold used (Q-stat  
292 p-value < 0.05). This indicates the complexity of these phenotypes, suggesting  
293 that multiple underlying mechanisms are likely involved. Evaluating multiple re-  
294 lated/correlated phenotypes is acceptable, as none of them are perfect measures of  
295 adiposity, and they likely proxy adiposity in different ways, e.g. BMI (subcuta-  
296 neous fat) vs waist-hip ratio (central adiposity).

297 Among physical activity traits (UK Biobank traits based on a survey question  
298 "*Types of physical activity in last 4 weeks*"), '*Heavy DIY*' shows evidence of a  
299 protective effect in overall samples (with little indication for heterogeneity and  
300 pleiotropy). Among dietary traits, there is a protective effect from '*dried fruit in-*  
301 *take*' on outcomes with larger sample sizes. The intake of *processed/sugary break-*  
302 *fast cereals* (i.e. *Frosties*, *Cornflakes*) appears to have a risk-increasing effect, but  
303 the instruments here are also likely highly pleiotropic as is commonly seen for di-  
304 etary intake traits. The important considerations for dietary and physical activity  
305 exposures in MR analyses are addressed in Appendix C.

306 The '*never smoked*' phenotype appears to confer a decreased risk in selected  
307 breast cancer subtypes. '*Chronotype (evening preference)*' shows evidence of in-  
308 creasing the risk in outcomes with larger sample sizes. '*Snoring*' also has a causal  
309 effect on some outcomes. Higher *age at menopause* has a risk-increasing effect on  
310 most outcomes, while '*Had menopause*' phenotype (potentially capturing earlier  
311 menopause) is protective.

### 312 3.1.2. Molecular traits

313 The lipid traits have a distinct effect direction pattern (Figure 4b): a consis-  
314 tently positive (causal) effect from HDL particle measures and a negative ( protec-  
315 tive) effect from VLDL. For all, there is limited evidence of effect on ER- breast

316 cancer and several molecular subtypes. It should be noted that LDL measures  
317 were also available in the MR-EvE discovery data, (i.e. 2218 traits), but did not  
318 pass filtering thresholds in the validation step due to weak evidence of effect (data  
319 available in Supplementary Table 3). Among other metabolites (Figure 4b), 3-  
320 dehydrocarnitine, betaine, and X-11440 show the strongest evidence of a negative  
321 effect on breast cancer risk.

322 The proteins with the most robust evidence (FDR-adjusted p-value < 0.05 and  
323 based on *cis*-instruments only) are IL1RL1, ISLR2, SAT2 (negative effect) and  
324 PRTN3, B3GNT2, CPXM1, PLXNB2 (positive effect). The proteins that did not  
325 pass FDR correction, but are also based on *cis*-only instruments and have effect  
326 on more than one outcome, include IL6R, EREG, DCBLD2 (negative effect) and  
327 VTN, IDUA, DNAJB11, XXYLT1 (positive effect).

328 Among the proteins for which genome-wide instruments were used (which  
329 may be picking up horizontal pleiotropy), the most consistent evidence of effect  
330 was observed for FLNA, ADH6, CAMK1 (negative effect) and IL21, CAST (posi-  
331 tive effect) with little heterogeneity and pleiotropy detected ('P' and 'H' not present  
332 in Figure 4c).

### 333 3.2. Evidence integration case studies

334 This section presents two exposure traits from the discovery of the MR-EvE  
335 risk factor, *childhood body size* and *HDL cholesterol*, as case studies. Although  
336 both have previously been investigated in relation to breast cancer, the underlying  
337 mechanisms remain unclear. We used MR-EvE data and data mined from liter-  
338 ature to generate hypotheses for potential mediators and cross-validate evidence  
339 from both sources. Since these traits have opposite effects on breast cancer risk  
340 and represent both molecular and non-molecular risk factors, they provide a ro-  
341 bust framework for demonstrating open and closed literature-based discovery, sup-  
342 ported by large literature spaces (Appendix A).

#### 343 3.2.1. Case study 1: Childhood body size

344 Childhood adiposity has been shown to have a protective effect on breast cancer  
345 risk in both observational [46, 47] and MR studies [48, 44], but the mechanism and  
346 mediators of this effect remain unclear. MR analysis shows strong evidence of an  
347 inverse effect from childhood body size (OpenGWAS ID: *ukb-b-4650*, female-only  
348 sample) on all breast cancer subtypes (Figure 4) (e.g. OR 0.60 [95% CI 0.53: 0.69]  
349 for BCAC 2017 overall sample).

350 **MR-EvE mediators.** In the initial MR-EvE scan, 125 putative mediators of the  
351 childhood body size effect on breast cancer were identified. Among these, several

352 were considered as potential mediators in earlier work [44]: IGF-1, age at menar-  
 353 che, age at menopause, age at last birth, glucose, glycated haemoglobin (Supple-  
 354 mentary Table 8). Following two-step MR results validation, 71 traits had evidence  
 355 of effect in both steps at nominal significance and 54 traits at the FDR-corrected  
 356 level (Supplementary Tables 9 and 10). Among these, a large proportion were adult  
 357 anthropometric traits (N=30). Since adult adiposity has been shown to be unlikely  
 358 to mediate the childhood adiposity effect [48, 49], these traits were not reviewed  
 359 further. Lipid traits (N=20) were also not considered further as they are unlikely to  
 360 mediate substantial effect individually. The remaining candidate mediators traits  
 361 were: *exercise to keep fit* (from UK Biobank physical activity questionnaire), *sleep*  
 362 *duration*, *skimmed milk intake*, and *Apolipoprotein A (ApoA)* (Table 1).

mediator	outcome	step1: beta	step1: p-value	step2: OR	step2: p-value
<b>Apolipoprotein A</b> ukb-d-30630_irtnt	Overall BC	-0.23 [-0.28:-0.17]	2.9e-15	0.77 [0.66:0.88]	2.2e-4
<b>Skimmed milk</b> ukb-d-1418_3	ER+ BC	0.09 [0.08:0.11]	2.1e-35	0.01 [0:0.15]	0.0014
<b>Exercise to keep fit</b> ukb-a-484	Overall BC	0.02 [0.001:0.04]	0.029	0.4 [0.21:0.75]	0.0044
<b>Sleep duration</b> ukb-b-4424	ER+ BC	-0.04 [-0.08:-0.01]	0.0053	1.49 [1.12:2.0]	0.0071

Table 1: **Two step-MR results for potential mediators of childhood body size/breast cancer relationship, identified in a hypothesis-free way via mining MR-EvE.** The results in the table are from the MR IVW validation analysis, showing beta [95%CI] for *step 1* (the effect of childhood body size on mediator, measured as SD change in mediator per body size category change) and OR [95%CI] for *step 2* (the odds of breast cancer per SD higher mediator). The table only includes the results that passed FDR in both steps. The table does not include anthropometric and lipid traits (the full results are available in Supplementary Table 9).

363 Higher childhood body size shows some evidence of an increasing effect on the  
 364 amount of *exercise to keep fit* in adulthood (effect size: 0.02 [95% CI 0.001: 0.04],  
 365 step 1), which in turn has a protective effect on breast cancer risk (OR: 0.4 [95%  
 366 CI 0.21: 0.75], step 2) (in agreement with published MR study [50] and results in  
 367 3.1.1). The relationship between childhood adiposity and the amount of exercise in  
 368 adulthood has not been explored with MR designs before, so this finding warrants  
 369 further follow-up.

370 Higher childhood body size has a negative effect on *sleep duration* in adult-  
 371 hood (effect size -0.04 [95% CI -0.08: -0.01], step 1), while longer sleeping time  
 372 increases breast cancer risk (stronger evidence for ER+, OR 1.49 [95% CI 1.12:  
 373 2.0], step 2), in agreement with published work [51].

374 **Literature mediators.** Next, we reviewed literature-mined terms that connect child-  
 375 hood adiposity to breast cancer. The Sankey diagram in Supplementary Figure B.6  
 376 shows the literature overlap between the two traits (Supplementary Table 17), per-  
 377 formed using the open discovery approach (Figure 3c, 2.6.2) (view interactively at



378 [https://mvab.shinyapps.io/literature\\_overlap\\_sankey/](https://mvab.shinyapps.io/literature_overlap_sankey/)).

379 The intermediate terms of the literature overlap include traits considered in  
380 earlier hypothesis-driven work [44] as potential mediators: IGF-1, testosterone,  
381 SHBG, insulin, and oestradiol. IGF-1 was also found using MR-EvE, but did not  
382 pass the FDR correction (mixed-sex sample data). The overlap contains many other  
383 literature-mined terms that make interesting mediation hypotheses (e.g. leptin, pro-  
384 lactin, adiponectin).

385 Next, to assess the identified intermediate terms as potential mediators using  
386 MR analysis, we searched OpenGWAS for suitable GWAS data to represent them.  
387 Of the 112 potential intermediate traits identified, data was available for 49, with  
388 38 traits having at least one SNP instrument (Supplementary Table 18).

389 We carried out MR bidirectionally between childhood body size (*ukb-b-4650*)  
390 and 38 potential intermediate traits (Supplementary Tables 19 and 20), followed  
391 by assessing their effect on breast cancer (Supplementary Table 21). A subset of  
392 combined results from the three MR analyses are presented in Table 2.

393 The criteria for a trait to have suggestive evidence of being a potential medi-  
394 ator (i.e. evidence of effect in *step1 forward* and *step2* (i.e. two-step MR), and  
395 preferably no effect in *step1 reverse*) was met by two traits: IGF-1 and bioavail-  
396 able testosterone. These two traits have been previously evaluated in [44] two-step  
397 MR analysis, with equivalent results and sensitivity analyses.

### 398 3.2.2. Case study 2: HDL-cholesterol

399 HDL-cholesterol (HDL-C) has been shown to have a risk-increasing effect on  
400 breast cancer in several MR studies [52, 53, 54], which opposes the finding from  
401 a large meta-analysis of prospective studies that showed a modestly inverse as-  
402 sociation [55]. This makes HDL-C an interesting case study for hypothesis-free  
403 mediators search, as this may help to dissect the mechanism linking it to breast  
404 cancer.

405 The HDL-C trait in MR-EvE is based on mixed-sex sample data from UK  
406 Biobank (OpenGWAS ID: *ukb-d-30760.irnt*); the evidence of effect from IVW  
407 for BCAC 2017 overall sample is OR 1.09 [95% 1.04: 1.13] (no effect on certain  
408 subtypes) (Figure 4b), which is in agreement with previous studies and consistent  
409 between sensitivity analysis MR methods (Supplementary Tables 4 and 6).

410 **MR-EvE mediators.** In the initial MR-EvE scan, 102 potential mediators of HDL-  
411 C's effect on breast cancer were identified, primarily lipid traits, proteins, and an-  
412 thropometric traits (Supplementary Table 11). After two-step MR validation, 60  
413 traits showed nominal significance, and 44 remained significant after FDR correc-  
414 tion (Supplementary Tables 12 and 13). Of these potential mediators, 36 were lipid  
415 traits, which we did not evaluate further due to their high correlation with HDL-C

Mediator trait	Step 1, forward (effect size, [95% CI])	Step 1, reverse (effect size,[95% CI])	Step 2 (OR, [95% CI])
<b>Alanine transaminase</b> ebi-a-GCST90013405	0.03 [0.02:0.04]*		
<b>CRP</b> ieu-b-35	0.39 [0.31:0.46]*		
<b>CSNK1D</b> prot-a-693			1.11 [1.01:1.22]
<b>Cortisone</b> met-a-354	-0.04 [-0.07:-0.02]*		
<b>E selectin</b> prot-b-77			0.95 [0.93:0.97]*
<b>EGF</b> prot-a-908		0.01 [0:0.02]	
<b>FABP4</b> prot-b-71	0.59 [0.27:0.90]*		
<b>Glycoprotein acetyls</b> met-c-863	0.33 [0.20:0.45]*		
<b>IGF-1</b> ukb-d-30770_irmt	-0.24 [-0.31:-0.16]*	-0.02 [-0.03:0]	1.07 [1.01:1.13]
<b>Insulin receptor</b> prot-a-1564			0.94 [0.92:0.97]*
<b>IGF1R</b> prot-a-1444			0.82 [0.76:0.88]*
<b>Leptin receptor</b> prot-a-1724		0.01 [0:0.01]	
<b>Oestradiol (F)</b> ieu-b-4872	-0.07 [-0.14:-0.01]	0.14 [0.08:0.21]*	
<b>SHBG (F)</b> ieu-b-4870	-0.20 [-0.29:-0.10]*		
<b>TSHbeta levels</b> ebi-a-GCST90010353		0.04 [0.01:0.08]	
<b>TNFRSF11B</b> prot-b-83			0.93 [0.87:0.99]
<b>Bio Testosterone (F)</b> ieu-b-4869	0.10 [0.04:0.17]*		1.12 [1.04:1.20]*
<b>Total Testosterone (F)</b> ieu-b-4864			1.14 [1.06:1.23]*

Table 2: **Subset of MR results (bidirectional and two-step) for potential mediators of childhood body size/breast cancer relationship, identified using LBD.** The table combines the results from three MR analyses: *step 1 forward* (childhood body size → mediator trait), *step 1 reverse* (mediator trait → childhood body size), and *step 2* (mediator trait → overall breast cancer). ‘*step 1 forward*’ and ‘*step 1 reverse*’ corresponded to bidirectional MR analysis; ‘*step 1 forward*’ and ‘*step 2*’ together correspond to two-step MR study design. The table only shows the results that pass the threshold of p.value < 0.05; the results that pass FDR are marked with an asterisk\*. The full results are available in Supplementary Tables 19, 20, 21. *Abbreviations:* CRP - C-Reactive protein, TSHbeta - Thyrotropin subunit beta, TNFRSF11b – tumour necrosis factor receptor superfamily member, CSNK1D - Casein kinase I isoform delta, EGF – epidermal growth factor, FABP4 - fatty acid binding protein 4; (F) - female-only data

416 and, therefore, the need for multivariable MR (MVMR) analysis. The identified  
 417 mediators also included several anthropometric traits, which may have a role in the  
 418 effect of HDL-C on breast cancer risk, but also would require further investigation  
 419 using MVMR for traits capturing different adiposity components (e.g., waist-hip  
 420 ratio, BMI).

421 HDL-C was found to have an increasing effect on two protein traits, *ApoA* (ef-

fect size: 0.87 [95% 0.84: 0.90]) and *cathepsin F* (CTSF) (0.14 [95% 0.03: 0.24]) (Table 1. The relationship and interplay between HDL-C, ApoA, and cathepsins have been investigated before [56, 57]. ApoA is the major protein component of HDL, which has an inverse effect on breast cancer in MR (OR: 0.77 [0.66:0.88], using *cis*-instruments only), also previously studied in connection to breast cancer [58]. CTSF has a positive effect on breast cancer risk (OR: 1.14 [95% 1.08: 1.21], single *cis*-instrument). CTSF belongs to a class of evolutionally conserved cysteine proteases involved in various biological processes that have been linked to disease, including cancer progression [59, 60]. Cathepsins (F, S, K; each to various degree) degrade lipid-free molecules of ApoA, leading to a complete loss of the ApoA's function as a cholesterol acceptor, stopping its anti-inflammatory action [56], which may be related to CTSF's role in disease. Further investigation into the mechanism linking HDL-C, ApoA, and CTSF (and potentially other cysteine proteases), and breast cancer is needed.

mediator	outcome	step 1: beta	step1: p-value	step2: OR	step2: p-value
<b>Apolipoprotein A</b> ukb-d-30630_irt	Overall BC	0.87 [0.84:0.9]	0	0.77 [0.66:0.88]	2.2e-4
<b>Cathepsin F (CTSF)</b> prot-a-722	Overall BC	0.14 [0.03:0.24]	0.014	1.14 [1.08:1.21]	7.5e-6
<b>Childhood obesity</b> ieu-a-1096	Overall BC	-0.19 [-0.35:-0.03]	0.023	0.82 [0.76:0.88]	2.4e-7
<b>Hip circumference</b> ieu-a-51	Overall BC	-0.07 [-0.11:-0.02]	0.0043	0.73 [0.59:0.9]	0.0039
<b>Waist circumference</b> ieu-a-63	Overall BC	-0.07 [-0.11:-0.02]	0.0039	0.73 [0.59:0.9]	0.0039
<b>Weight</b> ukb-b-12039	ER- BC	-0.04 [-0.07:-0.01]	0.0062	0.87 [0.79:0.95]	0.0029
<b>Exercise to keep fit</b> ukb-a-484	Overall BC	0.01 [0:0.02]	0.012	0.4 [0.21:0.75]	0.0044
<b>Snoring</b> ukb-a-14	ER- BC	0.01 [0:0.02]	0.0015	5.83 [1.9:17.87]	0.0021

Table 3: **Two step-MR results for potential mediators of HDL-cholesterol/breast cancer relationship, identified in a hypothesis-free way via mining MR-EvE.** The results in the table are from the MR IVW validation analysis, showing beta [95%CI] for *step 1* (the effect of HDL-C on mediator, measured as SD change in mediator per SD higher HDL-C) and OR [95%CI] for *step 2* (the odds of breast cancer per SD higher HDL-C). The table only includes the results that passed FDR in both steps. The table does not include lipid traits (the full results are available in Supplementary Table 12).

**Literature mediators.** Next, we reviewed literature-mined data connecting HDL and breast cancer to look for potential intermediates. Supplementary Figure B.7 shows the Sankey diagram of literature overlap between the traits (Supplementary Table 22), performed using the closed discovery approach (Figure 3b, 2.6.2) (view interactively at [https://mvab.shinyapps.io/literature\\_overlap\\_sankey/](https://mvab.shinyapps.io/literature_overlap_sankey/)).

The intermediate terms of this literature overlap include 120 unique (mostly

443 molecular) traits. Many of them are lipids-related traits (e.g., LDL, VLDL, ApoA,  
444 ApoB, Lipase); others include interesting hypotheses such as testosterone, IGF-1,  
445 insulin, CRP, and cathepsin D. Suitable GWAS datasets for MR validation were  
446 found for 50 of the 120 traits (Supplementary Table 23).

447 We carried out MR bidirectionally between HDL-C (*ukb-d-30760\_irnt*) and  
448 42 potential intermediate traits ( $\geq 1$  SNP) (Supplementary Tables 24 and 25),  
449 followed by assessing their effect on breast cancer (Supplementary Table 26). A  
450 subset of combined results from these three MR analyses is presented in Table 4.

451 From these results, the traits that fit the criteria for suggestive evidence of medi-  
452 ation are ApoA and testosterone. ApoA has been identified as a potential mediator  
453 using MR-EvE as well. HDL-C has opposite effects on two types of testosterone:  
454 negative on bioavailable and positive on total (Table 4), potentially indicating dif-  
455 ferent interaction mechanisms. Both testosterone measures have a risk-increasing  
456 effect on breast cancer (previously reported in [61, 44]). In the literature overlap,  
457 HDL is linked to testosterone via ApoE, HSD11B1, lipoproteins, and directly. In  
458 turn, testosterone is linked to breast cancer directly or via estrogen/progestin (Fig-  
459 ure B.7).

460 Finally, cathepsin D (CTSD) has been identified as a potential intermediate in  
461 the literature overlap (via CETP/endopeptidases and lipoprotein/cyclosporine, Fig-  
462 ure B.7). MR results show some evidence of HDL-C effect on CTSD (effect size  
463  $-0.12$  [95% CI  $-0.23$ : $-0.003$ ], but little evidence of CTSD effect on breast cancer  
464 (OR:  $1.02$  [95% CI  $0.99$ : $1.05$ ], three *cis*-SNPs) (Table 4, Supplementary Table 26),  
465 suggesting an unlikely role of CTSD as a mediator. However, the discovery of  
466 a CTSD in literature overlap complements the finding of CTSD in the MR-EvE  
467 data (Table 3), indicating that the cathepsins may have a role in the relationship of  
468 HDL-C and breast cancer, and the entire class of these enzymes should be explored  
469 in future work.

#### 470 **4. Discussion**

471 In this study, we used the EpiGraphDB knowledge graph to systematically  
472 identify breast cancer risk factors. Additionally, we integrated literature-mined  
473 and genetic data to determine potential mediators of the identified causal relation-  
474 ships between these risk factors and breast cancer, which we presented via two case  
475 studies.

476 Through mining MR-EvE data in EpiGraphDB we identified 2218 traits af-  
477 fecting breast cancer risk, which were narrowed down to 129 traits with reli-  
478 able causal evidence following rigorous validation (Figure 4, [https://mvab.shinyapps.io/MR\\_heatmaps/](https://mvab.shinyapps.io/MR_heatmaps/)). The identified risk/protective factors in-  
479 cluded both molecular biomarkers and non-molecular traits. Many of these traits  
480

Mediator trait	Step 1, forward (effect size,[95% CI])	Step 1, reverse (effect size, [95% CI])	Step 2 (OR, [95% CI])
<b>Adiponectin</b> ieu-a-1	0.1 [0.05:0.14]*		
<b>Albumin</b> met-c-841		0.05 [0:0.1]	
<b>Apolipoprotein A-I</b> met-c-842	0.71 [0.63:0.78]*	0.69 [0.60:0.77]*	0.64 [0.50:0.80]*
<b>Apolipoprotein B</b> ieu-b-108	-0.21 [-0.29:-0.12]*	-0.16 [-0.28:-0.05]*	
<b>Basigin</b> prot-a-275		0.11 [0.08:0.14]*	
<b>CRP</b> ukb-d-30710_irmt		0.03 [0.001:0.07]	
<b>CAMK1</b> prot-a-347		-0.03 [-0.05:-0.01]*	
<b>CAMK1</b> prot-a-346			0.98 [0.97:0.99]*
<b>Cathepsin D</b> prot-b-51	-0.12 [-0.23:-0.003]		
<b>Ephrin-B2</b> prot-a-905		-0.04 [-0.07:-0.01]	
<b>Histone deacetylase 8</b> prot-a-1320		-0.24 [-0.27:-0.21]*	
<b>IGF-1</b> ukb-d-30770_irmt			1.07 [1.01:1.13]*
<b>IGF1R</b> prot-a-1444		-0.1 [-0.13:-0.08]*	0.82 [0.76:0.88]*
<b>Lipoprotein A</b> ukb-d-30790_irmt		0.01 [0.001:0.02]*	
<b>Lipoprotein lipase</b> ebi-a-GCST90010149	0.24 [0.04:0.45]	0.27 [0.24:0.29]*	
<b>Oestradiol (F)</b> ieu-b-4872		0.28 [0.17:0.39]*	
<b>STARD1</b> prot-a-2866	-0.13 [-0.24:-0.01]		
<b>Superoxide dismutase</b> prot-a-2799	0.12 [0.01:0.24]	-0.01 [-0.02:-0.01]*	
<b>Bio Testosterone (F)</b> ieu-b-4869	-0.07 [-0.11:-0.03]*	-0.11 [-0.17:-0.05]*	1.12 [1.04:1.20]*
<b>Total Testosterone (F)</b> ieu-b-4864	0.03 [0.001:0.05]		1.14 [1.06:1.23]
<b>Total esterified cholesterol</b> met-d-Total_C_E	0.33 [0.27:0.39]*	0.30 [0.10:0.50]*	
<b>VLDL cholesterol</b> met-d-VLDL_C	-0.31 [-0.37:-0.24]*	-0.36 [-0.56:-0.17]*	

Table 4: **Subset of MR results (bidirectional and two-step) for potential mediators of HDL-C/breast cancer relationship, identified using LBD.** The table combines the results from three MR analyses: *step 1 forward* (HDL-C → mediator trait), *step 1 reverse* (mediator trait → HDL-C), and *step 2* (mediator trait → overall breast cancer). ‘*step 1 forward*’ and ‘*step 1 reverse*’ corresponded to bidirectional MR analysis; ‘*step 1 forward*’ and ‘*step 2*’ together correspond to two-step MR study design. The table only shows the results that pass the threshold of p.value < 0.05; the results that pass FDR are marked with an asterisk\*. The full results are available in Supplementary Tables 24,25,26. *Abbreviations:* CRP - C-Reactive protein, CAMK1 - Calcium/calmodulin-dependent protein kinase type 1, STARD1- Steroidogenic acute regulatory protein; (F) - female-only data

481 have been reported before, e.g. anthropometric measures [48, 62, 63], reproduc-  
482 tive traits [44], sleep traits [51, 64], height [65], smoking behaviour [66, 67],  
483 milk intake [68, 69], physical activity [70, 50, 71], lipids [52, 53, 54], metabo-  
484 lites [72], groups of proteins, e.g. adipokines [73], cytokines [74], inflammatory  
485 markers [75], metalloproteinases [76] and other circulating proteins [77, 78]; ad-  
486 ditionally, several studies evaluated multiple risk factor groups [44, 79, 80, 71].  
487 Several molecular traits that have not previously been reported in MR studies in-  
488 clude betaine, SAT2, PRTN3, PLXNB2, DCBLD2, IDUA. Appendix C provides  
489 an extended discussion of MR results of dietary and physical traits.

490 We presented two case study traits (*childhood body size* and *HDL-C*), for which  
491 we identified potential mediators from MR-EvE data and literature data using our  
492 LBD method, aiming to triangulate the evidence from both sources. Two-step  
493 MR and sensitivity analyses were used to provide suggestive evidence of media-  
494 tion through the identified traits. In the *childhood body size* case study, MR-EvE  
495 and literature-based mediator searches detected several overlapping traits, many of  
496 which have been considered as potential mediators in an earlier hypothesis-driven  
497 study [44] (e.g. IGF-1 and bioavailable testosterone), offering proof-of-concept  
498 evidence for this approach, in addition to identifying several traits not considered  
499 previously (e.g. physical activity and sleep duration). For *HDL-C*, the literature  
500 overlap method identified several intermediates that also appeared in the MR-EvE  
501 data (e.g. other lipid traits, ApoA, a cathepsin). HDL-C was also found to differen-  
502 tially affect total and bioavailable testosterone, which needs further investigation.  
503 Additionally, the entire class of cathepsin enzymes needs to be assessed for their  
504 mediating role in future work.

505 In our study, we took advantage of the pre-generated MR-EvE estimates in  
506 EpiGraphDB through the R package (*epigraphdb-r*) [4], and demonstrated how it  
507 enables fast hypothesis generation and prioritisation of candidate traits for more  
508 targeted investigation. Originally, MR-EvE was developed to provide a causal map  
509 of the human phenome, useful for supporting or rejecting specific hypotheses [14].  
510 However, for most traits, patterns of horizontal pleiotropy and issues with invalid  
511 instruments were detected [14]. This highlighted the need for rigorous follow-  
512 up, including sensitivity analyses, biologically informed instrument selection, and  
513 triangulation with other sources of evidence to validate the associations [14, 81].  
514 Hence, we validated our discoveries from MR-EvE using best practices in MR,  
515 accompanied by sensitivity analyses, *cis*-instrument extraction for protein traits,  
516 and also correction for multiple testing.

517 The MR-EvE data comes with several other limitations. Firstly, it is restricted  
518 to the traits available in OpenGWAS as of 2021, which were analysed in the  
519 everything-vs-everything way and subsequently added to EpiGraphDB. Secondly,  
520 a large proportion of those GWAS traits were measured in mixed-sex samples,

521 which suggests that traits with female-specific effects that are relevant for breast  
522 cancer risk could have been overlooked in the analysis (e.g. IGF-1 [44], testos-  
523 terone [61]). Moreover, many known risk factors (from observational or previous  
524 MR studies) may not come up in the MR-EvE review for several reasons, (1) some  
525 traits are not heritable/instrumentable or not in OpenGWAS, (2) traits were mea-  
526 sured in another (larger) cohort, or with additional covariates (e.g. BMI), (3) the  
527 filtering and FDR correction steps excluded traits with weaker evidence, which  
528 may have been reported by smaller studies (e.g. LDL-C).

529 Many traits identified using MR-EvE in this work have been published as indi-  
530 vidual MR studies over the recent years, despite MR-EvE data being available via  
531 EpiGraphDB since 2021 and earlier as a standalone resource (generated in 2017 by  
532 Hemani *et al* [14]) [82]. This highlights MR-EvE's value in generating hypotheses  
533 and guiding analyses, while also stressing the redundancy in the growing number  
534 of MR studies, which often report only two-sample MR results without providing  
535 further evidence, or attempting to dissect the underlying biological mechanisms  
536 [83, 84].

537 In our study, we address the need for evidence triangulation in MR studies by  
538 incorporating literature-mined data to strengthen the evidence for identified rela-  
539 tionships and to generate hypotheses for underlying mechanisms. To achieve this,  
540 we developed the LBD method (Figure 3) that helps to identify mechanistic path-  
541 ways or mediators for causal associations by mining literature data in EpiGraphDB.  
542 The identified intermediate terms may then be tested using MR (given suitable  
543 GWAS data is available). In both presented case studies, a third of the identi-  
544 fied intermediates provided testable hypotheses for MR validation. In this work,  
545 we limited GWAS data search to OpenGWAS due to practical constraints, but fu-  
546 ture research could extend this search to other sources (e.g. GWAS catalog [85])  
547 and the newly released molecular datasets (e.g. proteomics data (Olink) [86] and  
548 lipids/metabolites data (Nightingale NMR) [87]).

549 LBD and MR can mutually enhance each other: LBD can complement MR  
550 causal inference by identifying potential intermediates, while MR validation of  
551 LBD findings helps address challenges in evaluating discoveries and benchmark-  
552 ing the performance of LBD methods [88, 18]. Without proper evaluation, it's  
553 difficult to determine the usefulness or practicality of a discovery, which hinders  
554 the adoption of LBD in biomedical research [89]. Reviewing LBD discoveries with  
555 MR (if the terms are instrumentable and data is available) offers a relatively quick  
556 strategy for validation of the feasibility of the identified relationships. Certainly,  
557 MR is not a definitive answer for assessing traits' causal associations or validating  
558 a potential mediation mechanism. However, applying it to LBD discoveries may  
559 help to prioritise the most plausible mechanisms, which could be taken further into  
560 more costly and time-consuming validation.

561 Despite the potential gains of incorporating LBD into the biomedical/ epidemi-  
562 ological research process, the limitations and pitfalls of working with literature-  
563 mined data should be noted. Even though literature space extraction from EpiGraphDB  
564 is a straightforward process, the underlying raw triples data is noisy and may be  
565 difficult to interpret. The extraction of semantic triples from text, mapping terms  
566 to UMLS and later to GWAS phenotypes, is challenging and could be imprecise,  
567 and therefore this data must be treated with caution. Additionally, it is impor-  
568 tant to note that the chains of intermediate terms (presented as Sankey diagrams)  
569 should not be interpreted as true biological pathways. They simply represent re-  
570 lationships between biological entities mentioned in publications, and therefore  
571 direct interactions cannot be assumed. Moreover, the extracted relationships may  
572 also be affected by publication bias (i.e. due to abstracts often emphasising the  
573 established knowledge). Therefore, any insights drawn from literature-mined data  
574 must be traced back to the original publications to verify the validity of extracted  
575 relationships before being used as evidence.

## 576 **5. Conclusion**

577 In this article, we presented a novel and efficient approach to undertake a sys-  
578 tematic investigation of breast cancer risk factors by mining EpiGraphDB, repre-  
579 senting the largest MR study of breast cancer so far. We also introduced a novel  
580 LBD method for identifying potential intermediates or mechanisms between traits  
581 of interest. Our data-driven approach highlighted that MR and LBD can be used  
582 as complementary methods, providing supporting evidence to each other, as well  
583 as guiding the generation of hypotheses and bringing evidence from observational  
584 and genetic sources together. However, further validation through sensitivity tests,  
585 MVMR, and mediation analyses (where appropriate) is essential, as these are im-  
586 portant steps in validation of causal relationships. The risk factors and mediators  
587 identified in this work can be followed up in future studies in more detail, seeking  
588 evidence from other analyses or sources, such as observational studies or lab-based  
589 validation. Lastly, the overall approach can be used to study the aetiology of other  
590 diseases.

## 591 **6. CRediT authorship contribution statement**

592 **Marina Vabistsevits** - Conceptualization, Formal analysis, Data curation, In-  
593 vestigation, Methodology, Validation, Visualization, Writing – original draft, Writ-  
594 ing – review and editing; **Yi Liu** - Supervision, Writing – review and editing,  
595 Methodology, Software; **Tim Robinson** - Supervision, Writing – review and edit-  
596 ing; **Ben Elsworth** - Conceptualization, Supervision, Methodology, Software; **Tom**



597 **R Gaunt** - Conceptualization, Supervision, Methodology, Software, Writing – re-  
598 view and editing, Methodology, Project administration, Funding acquisition

## 599 **7. Declaration of Competing Interest**

600 T.R.G receives funding from Biogen and GSK for unrelated research.

## 601 **8. Acknowledgments**

602 M.V. is supported by the University of Bristol Alumni Fund (Professor Sir Eric  
603 Thomas Scholarship). T.R. is supported by NIHR Development and Skills En-  
604 hancement Award (NIHR 302363) and has received grants to attend educational  
605 workshops from Daiichi-Sankyo and Amgen. M.V., T.R., Y.L., T.R.G, work in the  
606 Medical Research Council Integrative Epidemiology Unit at the University of Bris-  
607 toल supported by the Medical Research Council (MC\_UU\_00032/03). This work  
608 was also supported by a Cancer Research UK programme grant (the Integrative  
609 Cancer Epidemiology Programme) (CC18281/A29019). This study was also sup-  
610 ported by the NIHR Biomedical Research Centre at University Hospitals Bristol  
611 NHS Foundation Trust and the University of Bristol. The views expressed in this  
612 publication are those of the author(s) and not necessarily those of the NHS, the  
613 National Institute for Health Research or the Department of Health.

## 614 **Appendix A. Literature overlap method details**

### 615 *Appendix A.1. Literature-mined relationships in EpiGraphDB - extended version*

616 EpiGraphDB contains literature-mined relationships between pairs of terms  
617 (also known as ‘literature triples’) that were derived from the published literature.  
618 Figure A.5 summarises the tools/databases involved in making literature-mined  
619 data available in EpiGraphDB. Literature-mined relationships in EpiGraphDB were  
620 originally extracted by the MELODI-Presto tool [38]. The underlying literature  
621 data comes from SemMedDB (Semantic MEDLINE Database) [16], a well-established  
622 repository of literature-mined semantic triples, i.e. ‘*subject term 1 – predicate –*  
623 *object term 2*’, mined from titles and abstracts of nearly 30 million biomedical arti-  
624 cles in PubMed, using SemRep [39, 40]. SemRep is a natural language processing  
625 (NLP) system that parses biomedical text using linguistic principles and UMLS  
626 (Unified Medical Language System) domain knowledge [90] to extract semantic  
627 triples.

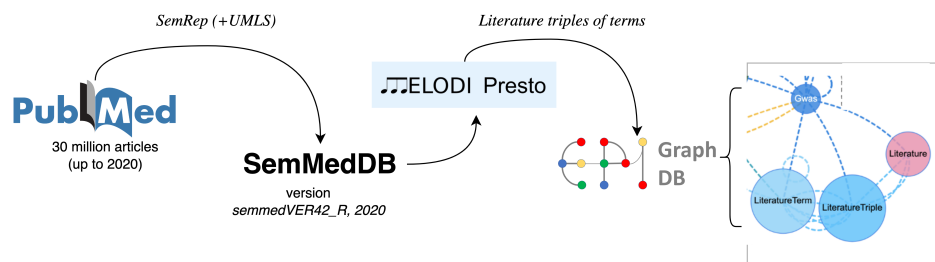


Figure A.5: A workflow diagram showing the origins of literature data in EpiGraphDB and the tools/databases involved in making this data available. EpiGraphDB literature term/triples nodes (on the right) is a subset of the full EpiGraphDB schema: <https://www.epigraphdb.org/about>

### 628 *Appendix A.2. Literature space overlap method - extended version*

629 Developing a method for connecting the literature spaces of exposure and out-  
630 come traits to explore potential links between them was one of the key aims of this  
631 work. The underlying principle of the literature space overlap approach is linking  
632 semantic triples of established biomedical knowledge into sequential ‘chains’, in  
633 order to uncover new connections. This principle of triple linking has previously  
634 been showcased in MELODI-Presto [38], and is an extension of the ABC paradigm  
635 [41, 43] with the expansion of a chain of Bs ( $B_1$  to  $B_i$ ):  $AB_1^iC$  (see 2.6.2). The  
636 details of the method are described below.

637 **Triple cleaning.** First, the extracted literature space for a trait is cleaned using  
638 several strategies – terms are standardised to minimise redundancy (i.e. merging  
639 long/short names, name acronyms, synonyms, gene/protein names); terms with  
640 different entity types (e.g. *aapp* - amino acid/protein, *ngm* - gene/genome) are  
641 filtered to include only the most common type; triples with the same *term1* and  
642 *term2* are excluded; among bidirectionally related terms the ones with more  
643 common direction (indicated by triple score) is retained. Additionally, terms that  
644 represented drugs/medications or other diseases (except anthropometrics-related  
645 terms) as well as negative predicates (e.g. ‘*neg coexists with*’) are excluded. Fi-  
646 nally, a ‘frequency score’ is generated for each pair of terms, i.e. the sum of triple  
647 scores across all possible predicates linking the terms, which is later used for triple  
648 chaining visualisation (Sankey diagrams).

649 **Term anchoring.** To implement *closed discovery* from the ABC paradigm [43]  
650 (see 2.6.2), we perform ‘term anchoring’, which involves restricting the outer term  
651 in a chain of triples to be the trait itself. For example, for an exposure trait (e.g.  
652 HDL-C), we are interested in downstream entities, so *term1* in the first triple of  
653 a chain make from its literature space would be ‘*HDL*’ (Figure 3b: “**T**” for trait).  
654 Inversely, in the breast cancer literature space, we are interested in entities upstream  
655 of breast cancer, so we would restrict *term2* of the terminal triple in the chain to be  
656 ‘*Breast cancer*’ (Figure 3b: “**BC**” for breast cancer).

657 For non-molecular traits, it is more difficult to perform term anchoring, because  
658 specific terms representing such traits are often not available. Therefore, for those  
659 we rely on *open discovery* [43] (see 2.6.2): the spaces are connected by match-  
660 ing unlinked triples from the exposure trait’s literature space to any breast cancer  
661 triples, and then extending chaining of the matched triples upstream into the trait’s  
662 space (Figure 3c).

663 **Triple chaining.** Literature spaces contain a variety of triples (some examples in  
664 Table A.5), and only a small proportion of them contain the trait itself as a term.  
665 Using only the triples that include the trait would make the analysis quite limited.  
666 Therefore, ‘triple chaining’ is implemented as a part of the approach to extend  
667 the discovery space. Within a trait’s literature space, we generate a basic set of  
668 triple chains, with the outer triple anchored to the trait term (see above), e.g. “**T**-  
669 A-B” and “C-D-E-F-**BC**”, where letters A-E are terms forming chains (see Figure  
670 3b). Then, the sets of triples from two literature spaces are “overlapped”, which  
671 involves *term2* from the exposure trait’s space triples to match with *term1* in triples  
672 from the breast cancer literature space. This can be done at any possible point, i.e.  
673 not necessarily as the longest multiple-step chain (e.g. “T-A-B-C-D-E-F-BC”), but  
674 can be of any variation (e.g. “T-A-B-BC”, “T-D-E-F-BC”).

term1	predicate	term2	triple score	PubMed ID (year)
estradiol	interacts_with	Estrogen Receptors	59	26466867 (2016) / 17975893 (2007) [...]
PTEN	interacts_with	Proto-Oncogene c-akt	21	23650412 (2013) / 27044817 (2016) [...]
estrogens	stimulates	Tamoxifen	18	10070945 (1999) / 28052658 (2017) [...]
Cyclin D1	interacts_with	estrogens	8	29935900 (2017) / 17638066 (2008) [...]
PGR	coexists_with	Estrogen Receptor alpha	4	29403307 (2018) / 26153859 (2015) [...]
CHEK1	interacts_with	BRCA1	2	31789405 (2019) / 12094328 (2002)
O-GlcNAc transferase	interacts_with	Cytoplasmic Protein	2	21255644 (2011) / 23576270 (2013)
polyglutamine	coexists_with	BRCA1	2	23483928 (2013) / 20422428 (2011)
Tretinoin	stimulates	Protein Kinase C	2	15516986 (2004) / 9284950 (1997)
Steroid Hormones	inhibits	SHBG	2	1320387 (1992) / 2173856 (1990)
AP2A1	stimulates	TCEAL1	1	16288208 (2006)
Protein Kinase C	stimulates	CTGF	1	16813525 (2006)
FOXMI	interacts_with	BIRC4	1	26404623 (2015)

Table A.5: **A subset of triples from the breast cancer literature space** (*term1* - *predicate* - *term2*) (randomly selected for illustrative purposes). The triple score indicates the number of publications the triple appears in/was extracted from. The 'PubMed ID (year)' column contains (a subset of) PubMed IDs of those publications. Full data is in Supplementary Table 14.

675 **Literature spaces.** The literature overlap method performs best on literature spaces  
 676 with  $> 50$  triples, as this is required for effective triple chaining. Small literature  
 677 spaces are also less likely to contain the trait as a term, which is required for term  
 678 anchoring. Therefore, literature space sizes were also considered when choosing  
 679 case study traits, to ensure an effective application of the literature overlap method.

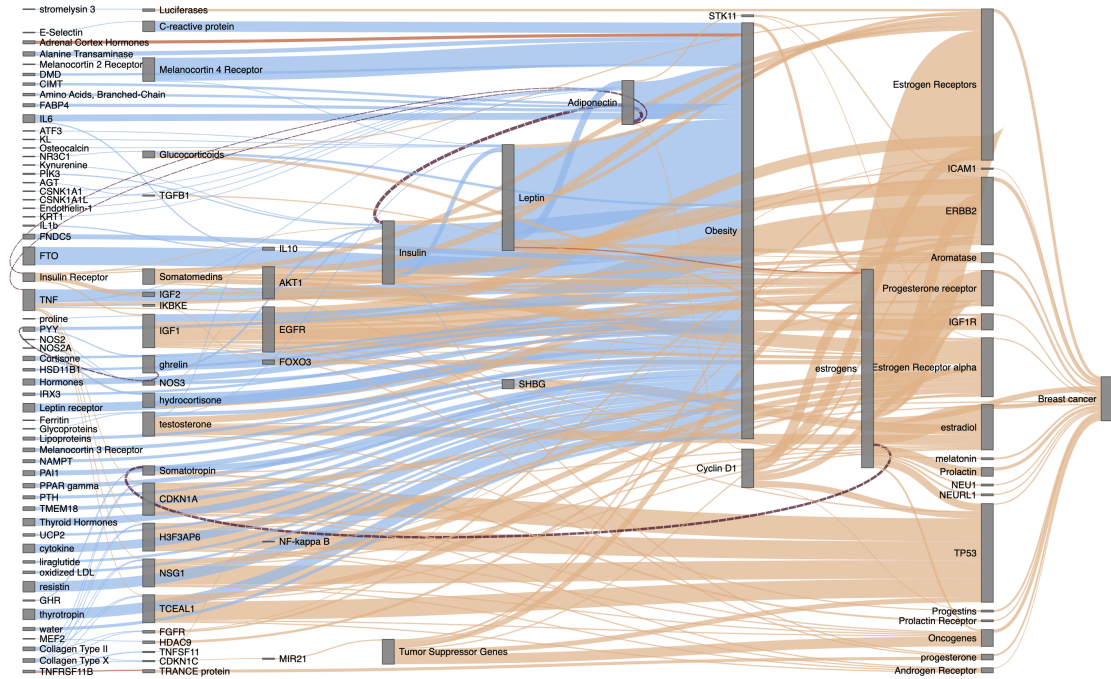
680 Literature spaces made of  $< 50$  triples indicate that fewer publications were  
 681 found to be linked to them. This may be due to the trait being less studied (and  
 682 published about), or issues with the trait name (i.e. non-standard phasing of a  
 683 concept), or issues with the source data in SemMedDB [16] or trait mapping in  
 684 MELODI-Presto [38] (Figure A.5). Literature spaces of this size are generally  
 685 unusable in the literature space overlap method.

### 686 Appendix A.3. Literature spaces - extended details

687 **Breast cancer.** Breast cancer literature space was formed from data linked to two  
 688 traits: 'Breast cancer' (OpenGWAS ID: *ieu-a-1126*) and 'Malignant neoplasm of  
 689 the breast' (*finn-a-C3\_BREAST*). 18,848 unique triples were identified, based on  
 690 4989 unique terms derived from 23,809 publications. The clean literature space  
 691 is available in Supplementary Table 14. The breast cancer anchor term is 'Breast  
 692 cancer', which is combined from the terms 'breast diseases' and 'malignant dis-  
 693 ease'.

694 **Case studies.** The selected case studies include *HDL* ('HDL-cholesterol' - *ukb-*  
 695 *d-30760\_irnt*, 7808 unique triples) and *childhood body size* ('childhood obesity'-  
 696 *ieu-a-1096*, 1255 unique triples) (Supplementary Tables 15-16). For childhood  
 697 body size, we chose to use 'childhood obesity' literature space over 'Comparative  
 698 body size at age 10' (*ukb-b-4650*, 211 triples) due to the larger space size and trait  
 699 name specificity.

700 **Appendix B. Sankey diagrams for case studies**



**Figure B.6: Sankey diagram of childhood obesity and breast cancer literature spaces overlap, showing the overlap of literature triples within and between the spaces.** The blue and orange relationships (triples) come from the childhood obesity and breast cancer literature spaces, respectively. The line width of each term relationship indicates how common it is in the literature (frequency score). Here, the open discovery approach was used, i.e. only one side of the overlap is anchored ('Breast cancer'); View interactively at [https://mvab.shinyapps.io/literature\\_overlap\\_sankey/](https://mvab.shinyapps.io/literature_overlap_sankey/)

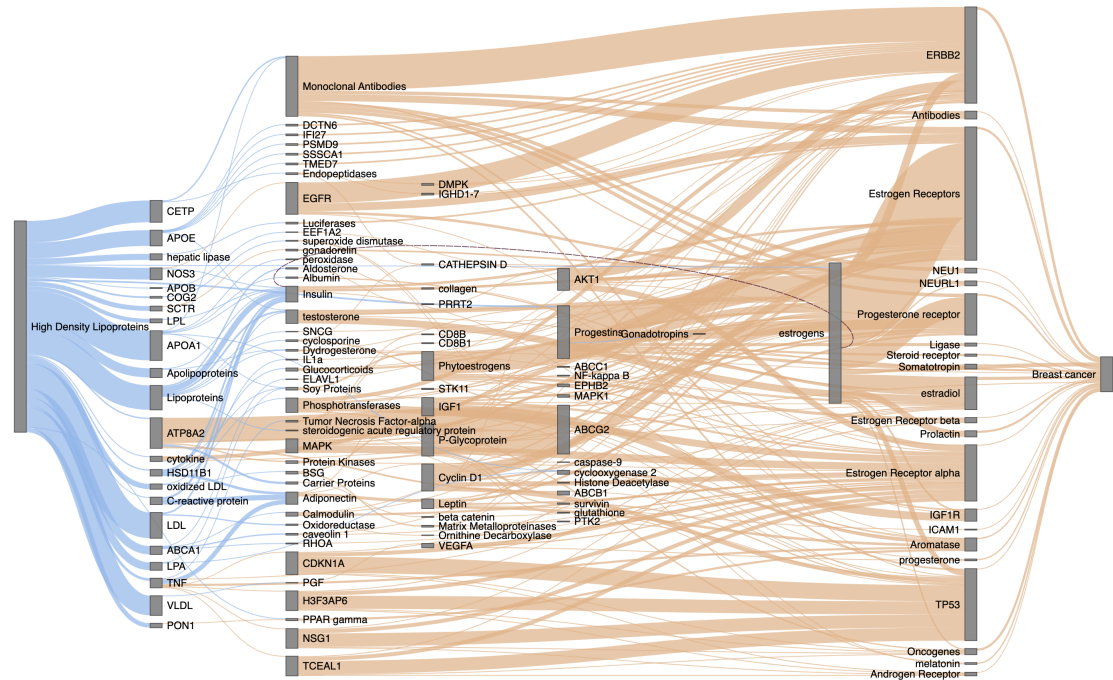


Figure B.7: Sankey diagram of HDL and breast cancer literature spaces overlap, showing how triples within and between HDL and breast cancer spaces interconnect. The blue and orange relationships (triples) come from the HDL and breast cancer literature spaces, respectively. The line width of each term relationship indicates how common it is in the literature (frequency score). Here, the closed discovery approach was used; the anchor terms are ‘High Density Lipoproteins’ and ‘Breast cancer’. View interactively at [https://mvab.shinyapps.io/literature\\_overlap\\_sankey/](https://mvab.shinyapps.io/literature_overlap_sankey/)

## 701 Appendix C. Extended discussion of MR results for dietary and physical traits

702 The final set of 129 validated breast cancer risk factors identified from MR-  
 703 EvE included several dietary traits. Genetically-proxied dietary intake exposures  
 704 require careful consideration due to the distinct sources of bias and interpretation  
 705 challenges they present. Dietary intake GWAS in UK Biobank can produce a con-  
 706 siderable number of instruments due to their sample size, and the risk of horizontal  
 707 pleiotropy among these SNPs is high [91]. Trait heterogeneity is also important  
 708 to consider [92], i.e. the fact that genetic instruments associated with the dietary  
 709 trait may represent different ‘dimensions’ of the same trait, e.g. intake of a certain  
 710 food item and metabolism of its core component (e.g. coffee/caffeine [93]). The  
 711 primary risk in over-interpreting dietary intake MR lies in the possibility that the

712 extracted instruments may be non-informative or associated with dietary patterns  
713 through socio-demographic and behavioural factors, potentially confounding the  
714 MR results [91]. Additionally, dietary trait instruments often include adiposity-  
715 related signals, meaning these variants could be affecting dietary patterns via their  
716 association with adiposity, leading to this association being captured in the MR  
717 results.

718 In this study, dietary traits with the most consistent evidence of effect were  
719 ‘*dried fruit intake*’ (showing a protective effect) and ‘*processed/sugary cereal in-*  
720 *take*’ (associated with an increased risk). Despite those MR results showing good  
721 instrument strength and no strong indication for pleiotropy but some heterogene-  
722 ity, we cannot claim that the captured effect is specific to those food items. The  
723 traits could be proxying two different environmental exposures, e.g. related to a  
724 broader class of foods that share similar characteristics, or overall dietary pattern,  
725 i.e. ‘healthy’ and ‘unhealthy’ food choices. The “gene-environment equivalence”  
726 assumption in MR states that the downstream effect of exposure modification is  
727 the same, regardless of it being genetically or non-genetically triggered. It could  
728 be hypothesised that the personal preference/choice to consume dried fruit vs sug-  
729 ary/processed cereal could be reflective of the other dietary patterns in those in-  
730 dividuals’ respective diets. Overall, the effect direction matches with the obser-  
731 vational evidence of diet impact on breast cancer risk [94]. The ‘healthy’ dietary  
732 pattern includes a varied diet with high fruit and vegetable intake, which increases  
733 the overall fibre and micronutrient content in the diet (more likely to occur in those  
734 consuming dried fruit), while the ‘unhealthy’ pattern involves a larger amount of  
735 processed and less nutrient-rich food items (more likely to occur in those choosing  
736 to eat sweetened breakfast cereals). The genetic correlation analysis in the study  
737 by Pirastu *et al* [95] has shown a positive correlation of ‘*dried fruit intake*’ GWAS  
738 with other traits relating to a healthier diet (e.g. intake of fresh fruit, cooked veg,  
739 fish, salad, being vegetarian) and a negative correlation with dietary patterns con-  
740 sidered less healthy (e.g. alcohol intake, processed meat, beef, spread on bread,  
741 salt, instant coffee), supporting the idea that ‘*dried fruit intake*’ instruments may  
742 be proxying an overall healthier diet. The protective effect of ‘*dried fruit intake*’ on  
743 breast cancer discovered from MR-EvE has been recently published in a separate  
744 MR study [96], in which the authors also reviewed the most likely confounders  
745 of this relationship (e.g. vitamin C, BMI, years of education) using MVMR, and  
746 found that the identified inverse relationship is not affected when adjusted for any  
747 of them.

748 *Skimmed milk intake* was another dietary trait identified from MR-EvE with  
749 evidence of a protective effect on breast cancer, but this result has to be carefully  
750 considered as it is misleading. Previous MR studies of milk intake effect [68, 69]  
751 selected instruments (nSNP=1) specifically within the lactase gene locus (*LCT*),

752 which encodes an enzyme that ensures tolerance of milk products. The variation  
753 in this locus affects the amount of milk consumed and, therefore, is often used as a  
754 proxy for this exposure. Both MR studies found a mild risk-increasing effect from  
755 milk intake. In MR-EvE data, the instruments for ‘*skimmed milk intake*’ trait were  
756 extracted genome-wide (nSNP=3), none of them being in *LCT*. MR sensitivity  
757 tests identified likely horizontal pleiotropy and heterogeneity. This serves as a  
758 valuable negative example of dietary traits evaluation and highlights the importance  
759 of conducting more detailed investigations of these traits in individual studies.

760 Physical activity is one of the modifiable lifestyle factors that have been shown  
761 to have a protective effect on breast cancer risk, particularly vigorous activity [97].  
762 MR studies using accelerometer-measured overall physical activity produced re-  
763 sults to support this [70, 50]. Among the physical activity traits available in MR-  
764 EvE (UK Biobank traits based on a survey question “*Types of physical activity*  
765 *in last 4 weeks*”), several showed similar protective effects: ‘*Heavy DIY*’, ‘*stren-*  
766 *uous sports*’ and ‘*exercise to keep fit*’ (the last identified in a case study), which  
767 may be proxying moderate-to-vigorous physical activity. Similarly, there was a  
768 risk-increasing effect from ‘*no physical activity*’ trait, which is in agreement with  
769 ‘*sedentary time*’ found to increase the risk in [50] and the observationally-known  
770 detrimental effects of physical inactivity [98]. With survey-derived physical activ-  
771 ity phenotypes (such as ‘*Heavy DIY*’) appearing unreliable, a potential follow-up  
772 to these findings would be to compare the genetic instruments identified for physi-  
773 cal activity from survey phenotypes to those derived from accelerometer-measured  
774 data.

## 775 References

- 776 [1] D. A. Lawlor, K. Tilling, G. D. Smith, Triangulation in aetiological epi-  
777 demiology, *International Journal of Epidemiology* 45 (6) (2017) 1866–1886.  
778 doi:10.1093/ije/dyw314.  
779 URL [https://academic.oup.com/ije/article/45/6/1866/](https://academic.oup.com/ije/article/45/6/1866/2930550)  
780 [2930550](https://academic.oup.com/ije/article/45/6/1866/2930550)
- 781 [2] M. R. Munafò, G. Davey Smith, Robust research needs many lines of  
782 evidence (1 2018). doi:10.1038/d41586-018-01023-3.  
783 URL [https://www.nature.com/articles/](https://www.nature.com/articles/d41586-018-01023-3)  
784 [d41586-018-01023-3](https://www.nature.com/articles/d41586-018-01023-3)
- 785 [3] M. R. Munafò, J. P. Higgins, G. D. Smith, Triangulating evidence through the  
786 inclusion of genetically informed designs, *Cold Spring Harbor Perspectives*  
787 *in Medicine* 11 (8) (2021). doi:10.1101/cshperspect.a040659.  
788 URL <https://pubmed.ncbi.nlm.nih.gov/33355252/>



- 789 [4] Y. Liu, B. Elsworth, P. Erola, V. Haberland, G. Hemani, M. Lyon, J. Zheng,  
790 O. Lloyd, M. Vabistsevits, T. R. Gaunt, EpiGraphDB: A database and data  
791 mining platform for health data science, *Bioinformatics* 37 (9) (2021) 1304–  
792 1311. doi:10.1093/bioinformatics/btaa961.  
793 URL <https://doi.org/10.1101/2020.08.01.230193>
- 794 [5] D. Ochoa, A. Hercules, M. Carmona, D. Suveges, A. Gonzalez-Uriarte,  
795 C. Malangone, A. Miranda, L. Fumis, D. Carvalho-Silva, M. Spitzer,  
796 J. Baker, J. Ferrer, A. Raies, O. Razuvayevskaya, A. Faulconbridge,  
797 E. Petsalaki, P. Mutowo, S. MacHlitt-Northen, G. Peat, E. McAuley, C. K.  
798 Ong, E. Mountjoy, M. Ghoussaini, A. Pierleoni, E. Papa, M. Pignatelli,  
799 G. Koscielny, M. Karim, J. Schwartzentruber, D. G. Hulcoop, I. Dunham,  
800 E. M. McDonagh, Open Targets Platform: Supporting systematic drug-target  
801 identification and prioritisation, *Nucleic Acids Research* 49 (D1) (2021)  
802 D1302–D1310. doi:10.1093/nar/gkaa1027.  
803 URL [https://academic.oup.com/nar/article/49/D1/  
804 D1302/5983621](https://academic.oup.com/nar/article/49/D1/D1302/5983621)
- 805 [6] A. Struck, B. Walsh, A. Buchanan, J. A. Lee, R. Spangler, J. M. Stu-  
806 art, K. Ellrott, Exploring Integrative Analysis Using the BioMedical  
807 Evidence Graph, *JCO Clinical Cancer Informatics* 4 (2020) 147–159.  
808 doi:10.1200/cci.19.00110.  
809 URL [https://ascopubs.org/doi/full/10.1200/CCI.19.  
810 00110](https://ascopubs.org/doi/full/10.1200/CCI.19.00110)
- 811 [7] E. A. Coker, C. Mitsopoulos, J. E. Tym, A. Komianou, C. Kannas,  
812 P. Di Micco, E. Villasclaras Fernandez, B. Ozer, A. A. Antolin, P. Work-  
813 man, B. Al-Lazikani, CanSAR: Update to the cancer translational research  
814 and drug discovery knowledgebase, *Nucleic Acids Research* 47 (D1) (2019)  
815 D917–D922. doi:10.1093/nar/gky1129.  
816 URL <https://pubmed.ncbi.nlm.nih.gov/33219674/>
- 817 [8] S. Ebrahim, G. Davey Smith, Mendelian randomization: Can genetic  
818 epidemiology help redress the failures of observational epidemi-  
819 ology?, *International Journal of Epidemiology* 32 (1) (2003) 1–22.  
820 doi:10.1007/s00439-007-0448-6.  
821 URL [https://link.springer.com/article/10.1007/  
822 s00439-007-0448-6](https://link.springer.com/article/10.1007/s00439-007-0448-6)
- 823 [9] E. Sanderson, M. M. Glymour, M. V. Holmes, H. Kang, J. Morrison, M. R.  
824 Munafò, T. Palmer, C. M. Schooling, C. Wallace, Q. Zhao, G. Davey Smith,  
825 Mendelian randomization, *Nature Reviews Methods Primers* 2 (1) (12 2022).

- 826 doi:10.1038/s43586-021-00092-5.  
827 URL <https://www.nature.com/articles/s43586-021-00092-5>  
828
- 829 [10] J. Xu, M. Li, Y. Gao, M. Liu, S. Shi, J. Shi, K. Yang, Z. Zhou, J. Tian,  
830 Using Mendelian randomization as the cornerstone for causal inference in  
831 epidemiology, *Environmental Science and Pollution Research* 29 (4) (2022)  
832 5827–5839. doi:10.1007/s11356-021-15939-3.  
833 URL <https://pubmed.ncbi.nlm.nih.gov/34431050/>
- 834 [11] B. Elsworth, M. Lyon, T. Alexander, Y. Liu, P. Matthews, J. Hallett, P. Bates,  
835 T. Palmer, V. Haberland, G. Davey Smith, J. Zheng, P. Haycock, T. R. Gaunt,  
836 G. Hemani, B. Elsworth, M. Lyon, T. Alexander, Y. Liu, P. Matthews, J. Hal-  
837 lett, P. Bates, T. Palmer, V. Haberland, G. D. Smith, J. Zheng, P. Haycock,  
838 T. R. Gaunt, G. Hemani, The MRC IEU OpenGWAS data infrastructure,  
839 *bioRxiv* (8 2020). doi:10.1101/2020.08.10.244293.  
840 URL <https://doi.org/10.1101/2020.08.10.244293>
- 841 [12] E. Sollis, A. Mosaku, A. Abid, A. Buniello, M. Cerezo, L. Gil, T. Groza,  
842 O. Güneş, P. Hall, J. Hayhurst, A. Ibrahim, Y. Ji, S. John, E. Lewis,  
843 J. A. Macarthur, A. McMahon, D. Osumi-Sutherland, K. Panoutsopoulou,  
844 Z. Pendlington, S. Ramachandran, R. Stefancsik, J. Stewart, P. Whetzel,  
845 R. Wilson, L. Hindorff, F. Cunningham, S. A. Lambert, M. Inouye, H. Parkin-  
846 son, L. W. Harris, The NHGRI-EBI GWAS Catalog: knowledgebase and  
847 deposition resource, *Nucleic Acids Research* 51 (D1) (2023) D977–D985.  
848 doi:10.1093/NAR/GKAC1010.  
849 URL <https://dx.doi.org/10.1093/nar/gkac1010>
- 850 [13] G. Hemani, J. Zheng, B. Elsworth, K. H. Wade, V. Haberland, D. Baird,  
851 C. Laurin, S. Burgess, J. Bowden, R. Langdon, V. Y. Tan, J. Yarmolin-  
852 sky, H. A. Shihab, N. J. Timpson, D. M. Evans, C. Relton, R. M. Martin,  
853 G. Davey Smith, T. R. Gaunt, P. C. Haycock, The MR-base platform supports  
854 systematic causal inference across the human phenome, *eLife* 7 (5 2018).  
855 doi:10.7554/eLife.34408.  
856 URL <https://pubmed.ncbi.nlm.nih.gov/29846171/>
- 857 [14] G. Hemani, J. Bowden, P. Haycock, J. Zheng, O. Davis, P. Flach, T. Gaunt,  
858 G. D. Smith, Automating Mendelian randomization through machine learn-  
859 ing to construct a putative causal map of the human phenome, *bioRxiv* (10  
860 (2017) 173682. doi:10.1101/173682.  
861 URL [https://www.biorxiv.org/content/10.1101/  
862 173682v2](https://www.biorxiv.org/content/10.1101/173682v2)

- 863 [15] L. J. Jensen, J. Saric, P. Bork, Literature mining for the biologist: From infor-  
864 mation retrieval to biological discovery (2006). doi:10.1038/nrg1768.  
865 URL [www.nature.com/reviews/genetics](http://www.nature.com/reviews/genetics)
- 866 [16] H. Kilicoglu, D. Shin, M. Fiszman, G. Roseblat, T. C. Rindfleisch,  
867 SemMedDB: A PubMed-scale repository of biomedical semantic predi-  
868 cations, *Bioinformatics* 28 (23) (2012) 3158–3160. doi:10.1093/  
869 bioinformatics/bts591.  
870 URL <https://pubmed.ncbi.nlm.nih.gov/23044550/>
- 871 [17] T. Bekhuis, Conceptual biology, hypothesis discovery, and text min-  
872 ing: Swanson’s legacy, *Biomedical Digital Libraries* 3 (1) (2006) 1–7.  
873 doi:10.1186/1742-5581-3-2/METRICS.  
874 URL [https://bio-diglib.biomedcentral.com/articles/](https://bio-diglib.biomedcentral.com/articles/10.1186/1742-5581-3-2)  
875 [10.1186/1742-5581-3-2](https://bio-diglib.biomedcentral.com/articles/10.1186/1742-5581-3-2)
- 876 [18] G. Crichton, S. Baker, Y. Guo, A. Korhonen, Neural networks for open and  
877 closed Literature-based Discovery, *PLoS ONE* 15 (5) (2020) e0232891.  
878 doi:10.1371/journal.pone.0232891.  
879 URL [https://journals.plos.org/plosone/article?id=](https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0232891)  
880 [10.1371/journal.pone.0232891](https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0232891)
- 881 [19] WCRF, World Cancer Research Fund/American Institute for Cancer Re-  
882 search. Continuous Update Project. Diet, Nutrition, Physical Activity and the  
883 Prevention of Cancer. Summary of evidence. (2018).  
884 URL [www.wcrf.org/matrix](http://www.wcrf.org/matrix)
- 885 [20] IARC, International Agency for Research on Cancer: Estimated cumulative  
886 risk of incidence in 2020, in females, in high-income countries, by cancer  
887 site; based on GLOBOSCAN 2020 data (2021).  
888 URL [https://gco.iarc.fr/today/](https://gco.iarc.fr/today/online-analysis-multi-bars?v=2020&mode=cancer&mode_population=countries&population=900&populations=986&key=cum_risk&sex=2&cancer=39&type=0&statistic=5&prevalence=0&population_group=0&ages_group%5B%5D=0&ages_group%5B%5D=14&nb_item)  
889 [online-analysis-multi-bars?v=2020&mode=cancer&](https://gco.iarc.fr/today/online-analysis-multi-bars?v=2020&mode=cancer&mode_population=countries&population=900&populations=986&key=cum_risk&sex=2&cancer=39&type=0&statistic=5&prevalence=0&population_group=0&ages_group%5B%5D=0&ages_group%5B%5D=14&nb_item)  
890 [mode\\_population=countries&population=900&](https://gco.iarc.fr/today/online-analysis-multi-bars?v=2020&mode=cancer&mode_population=countries&population=900&populations=986&key=cum_risk&sex=2&cancer=39&type=0&statistic=5&prevalence=0&population_group=0&ages_group%5B%5D=0&ages_group%5B%5D=14&nb_item)  
891 [populations=986&key=cum\\_risk&sex=2&cancer=39&type=](https://gco.iarc.fr/today/online-analysis-multi-bars?v=2020&mode=cancer&mode_population=countries&population=900&populations=986&key=cum_risk&sex=2&cancer=39&type=0&statistic=5&prevalence=0&population_group=0&ages_group%5B%5D=0&ages_group%5B%5D=14&nb_item)  
892 [0&statistic=5&prevalence=0&population\\_group=0&ages\\_](https://gco.iarc.fr/today/online-analysis-multi-bars?v=2020&mode=cancer&mode_population=countries&population=900&populations=986&key=cum_risk&sex=2&cancer=39&type=0&statistic=5&prevalence=0&population_group=0&ages_group%5B%5D=0&ages_group%5B%5D=14&nb_item)  
893 [group%5B%5D=0&ages\\_group%5B%5D=14&nb\\_item](https://gco.iarc.fr/today/online-analysis-multi-bars?v=2020&mode=cancer&mode_population=countries&population=900&populations=986&key=cum_risk&sex=2&cancer=39&type=0&statistic=5&prevalence=0&population_group=0&ages_group%5B%5D=0&ages_group%5B%5D=14&nb_item)
- 894 [21] K. L. Britt, J. Cuzick, K. A. Phillips, Key steps for effective breast  
895 cancer prevention, *Nature Reviews Cancer* (2020). doi:10.1038/  
896 s41568-020-0266-x.  
897 URL <http://dx.doi.org/10.1038/s41568-020-0266-x>

- 898 [22] N. Harbeck, F. Penault-Llorca, J. Cortes, M. Gnant, N. Houssami, P. Poort-  
899 mans, K. Ruddy, J. Tsang, F. Cardoso, Breast cancer, *Nature Reviews Disease*  
900 *Primers* 5 (1) (2019). doi:10.1038/s41572-019-0111-2.  
901 URL <https://pubmed.ncbi.nlm.nih.gov/31548545/>
- 902 [23] C. L. Relton, G. Davey Smith, Two-step epigenetic mendelian randomiza-  
903 tion: A strategy for establishing the causal role of epigenetic processes in  
904 pathways to disease, *International Journal of Epidemiology* 41 (1) (2012)  
905 161–176. doi:10.1093/ije/dyr233.  
906 URL <https://pubmed.ncbi.nlm.nih.gov/22422451/>
- 907 [24] J. Zheng, D. Baird, M.-C. Borges, J. Bowden, G. Hemani, P. Haycock,  
908 D. M. Evans, G. D. Smith, Recent Developments in Mendelian Random-  
909 ization Studies, *Current Epidemiology Reports* 4 (4) (2017) 330–345. doi:  
910 10.1007/s40471-017-0128-6.  
911 URL <https://doi.org/10.1007/s40471-017-0128-6>
- 912 [25] K. Michailidou, S. Lindström, J. Dennis, J. Beesley, S. Hui, S. Kar,  
913 A. Lemaçon, P. Soucy, D. Glubb, A. Rostamianfar, M. K. Bolla, Q. Wang,  
914 J. Tyrer, E. Dicks, A. Lee, Z. Wang, J. Allen, R. Keeman, U. Eilber, J. D.  
915 French, X. Qing Chen, L. Fachal, K. McCue, A. E. McCart Reed, M. Ghous-  
916 saine, J. S. Carroll, X. Jiang, H. Finucane, M. Adams, M. A. Adank, H. Ahsan,  
917 K. Aittomäki, H. Anton-Culver, N. N. Antonenkova, V. Arndt, K. J. Aronson,  
918 B. Arun, P. L. Auer, F. Bacot, M. Barrdahl, C. Baynes, M. W. Beckmann,  
919 S. Behrens, J. Benitez, M. Bermisheva, L. Bernstein, C. Blomqvist, N. V.  
920 Bogdanova, S. E. Bojesen, B. Bonanni, A.-L. Børresen-Dale, J. S. Brand,  
921 H. Brauch, P. Brennan, H. Brenner, L. Brinton, P. Broberg, I. W. Brock,  
922 A. Broeks, A. Brooks-Wilson, S. Y. Brucker, T. Brüning, B. Burwinkel,  
923 K. Butterbach, Q. Cai, H. Cai, T. Caldés, F. Canzian, A. Carracedo, B. D.  
924 Carter, J. E. Castelao, T. L. Chan, T.-Y. David Cheng, K. Seng Chia, J.-Y.  
925 Choi, H. Christiansen, C. L. Clarke, M. Collée, D. M. Conroy, E. Cordina-  
926 Duverger, S. Cornelissen, D. G. Cox, A. Cox, S. S. Cross, J. M. Cunning-  
927 ham, K. Czene, M. B. Daly, P. Devilee, K. F. Doheny, T. Dörk, I. dos  
928 Santos-Silva, M. Dumont, L. Durcan, M. Dwek, D. M. Eccles, A. B. Ekici,  
929 A. Heather Eliassen, C. Ellberg, M. Elvira, C. Engel, M. Eriksson, P. A.  
930 Fasching, J. Figueroa, D. Flesch-Janys, O. Fletcher, H. Flyger, L. Fritschi,  
931 V. Gaborieau, M. Gabrielson, M. Gago-Dominguez, Y.-T. Gao, S. M. Gap-  
932 stur, J. A. García-Sáenz, M. M. Gaudet, V. Georgoulas, G. G. Giles, G. Glend-  
933 on, M. S. Goldberg, D. E. Goldgar, A. González-Neira, G. I. Grenaker Al-  
934 næs, M. Grip, J. Gronwald, A. Grundy, P. Guénel, L. Haeberle, E. Hahn-  
935 en, C. A. Haiman, N. Håkansson, U. Hamann, N. Hamel, S. Hankin-

936 son, P. Harrington, S. N. Hart, J. M. Hartikainen, M. Hartman, A. Hein,  
937 J. Heyworth, B. Hicks, P. Hillemanns, D. N. Ho, A. Hollestelle, M. J.  
938 Hooning, R. N. Hoover, J. L. Hopper, M.-F. Hou, C.-N. Hsiung, G. Huang,  
939 K. Humphreys, J. Ishiguro, H. Ito, M. Iwasaki, H. Iwata, A. Jakubowska,  
940 W. Janni, E. M. John, N. Johnson, K. Jones, M. Jones, A. Jukkola-Vuorinen,  
941 R. Kaaks, M. Kabisch, K. Kaczmarek, D. Kang, Y. Kasuga, M. J. Kerin,  
942 S. Khan, E. Khusnutdinova, J. I. Kiiski, S.-W. Kim, J. A. Knight, V.-M.  
943 Kosma, V. N. Kristensen, U. Krüger, A. Kwong, D. Lambrechts, L. Le Marc-  
944 hand, E. Lee, M. Hyuk Lee, J. Won Lee, C. Neng Lee, F. Lejbkowitz,  
945 J. Li, J. Lilyquist, A. Lindblom, J. Lissowska, W.-Y. Lo, S. Loibl, J. Long,  
946 A. Lophatananon, J. Lubinski, C. Luccarini, M. P. Lux, E. S. K Ma, R. J.  
947 MacInnis, T. Maishman, E. Makalic, K. E. Malone, I. Maleva Kostovska,  
948 A. Mannermaa, S. Manoukian, J. E. Manson, S. Margolin, S. Mariapun,  
949 M. Elena Martinez, K. Matsuo, D. Mavroudis, J. McKay, C. McLean,  
950 H. Meijers-Heijboer, A. Meindl, P. Menéndez, U. Menon, J. Meyer, H. Miao,  
951 N. Miller, N. Aishah Mohd Taib, K. Muir, A. Marie Mulligan, C. Mulot,  
952 S. L. Neuhausen, H. Nevanlinna, P. Neven, S. F. Nielsen, D.-Y. Noh, B. G.  
953 Nordestgaard, A. Norman, O. I. Olopade, J. E. Olson, H. Olsson, C. Olswold,  
954 N. Orr, V. Shane Pankratz, S. K. Park, T.-W. Park-Simon, R. Lloyd, J. I.  
955 A Perez, P. Peterlongo, J. Peto, K.-A. Phillips, M. Pinchev, D. Plaseska-  
956 Karanfilska, R. Prentice, N. Presneau, D. Prokofyeva, E. Pugh, K. Pylkäs,  
957 B. Rack, P. Radice, N. Rahman, G. Rennert, H. S. Rennert, V. Rhenius,  
958 A. Romero, J. Romm, K. J. Ruddy, T. Rüdiger, A. Rudolph, M. Ruebner,  
959 E. J. T Rutgers, E. Saloustros, D. P. Sandler, S. Sangrajang, E. J. Sawyer,  
960 D. F. Schmidt, R. K. Schmutzler, A. Schneeweiss, M. J. Schoemaker, F. Schu-  
961 macher, P. Schürmann, R. J. Scott, C. Scott, S. Seal, C. Seynaeve, M. Shah,  
962 P. Sharma, C.-Y. Shen, G. Sheng, M. E. Sherman, M. J. Shrubsole, X.-O. Shu,  
963 A. Smeets, C. Sohn, M. C. Southey, J. J. Spinelli, C. Stegmaier, S. Stewart-  
964 Brown, J. Stone, D. O. Stram, H. Surowy, A. Swerdlow, R. Tamimi, J. A.  
965 Taylor, M. Tengström, S. H. Teo, M. Beth Terry, D. C. Tessier, S. Thana-  
966 sitthichai, K. Thöne, R. A. E M Tollenaar, I. Tomlinson, L. Tong, D. Tor-  
967 res, T. Truong, C.-C. Tseng, S. Tsugane, H.-U. Ulmer, G. Ursin, M. Untch,  
968 C. Vachon, C. J. van Asperen, D. Van Den Berg, A. M. W van den Ouweland,  
969 L. van der Kolk, R. B. van der Luit, D. Vincent, J. Vollenweider, Q. Waisfisz,  
970 S. Wang-Gohrke, C. R. Weinberg, C. Wendt, A. S. Whittemore, H. Wildiers,  
971 W. Willett, R. Winqvist, A. Wolk, A. H. Wu, L. Xia, T. Yamaji, X. R. Yang,  
972 C. Har Yip, K.-Y. Yoo, J.-C. Yu, W. Zheng, Y. Zheng, B. Zhu, A. Ziogas,  
973 E. Ziv, S. R. Lakhani, A. C. Antoniou, A. Droit, I. L. Andrulis, C. I. Amos,  
974 F. J. Couch, P. D. P Pharoah, J. Chang-Claude, P. Hall, D. J. Hunter, R. L.  
975 Milne, M. García-Closas, M. K. Schmidt, S. J. Chanock, A. M. Dunning,

- 976 S. L. Edwards, G. D. Bader, G. Chenevix-Trench, J. Simard, P. Kraft, D. F.  
977 Easton, Association analysis identifies 65 new breast cancer risk loci, *Nature*  
978 (2017). doi:10.1038/nature24284.  
979 URL [www.nature.com/reprints](http://www.nature.com/reprints) .
- 980 [26] H. Zhang, T. U. Ahearn, J. Lecarpentier, D. Barnes, J. Beesley, G. Qi,  
981 X. Jiang, T. A. O’Mara, N. Zhao, M. K. Bolla, A. M. Dunning, J. Dennis,  
982 Q. Wang, Z. A. Ful, K. Aittomäki, I. L. Andrulis, H. Anton-Culver, V. Arndt,  
983 K. J. Aronson, B. K. Arun, P. L. Auer, J. Azzollini, D. Barrowdale, H. Becher,  
984 M. W. Beckmann, S. Behrens, J. Benitez, M. Bermisheva, K. Bialkowska,  
985 A. Blanco, C. Blomqvist, N. V. Bogdanova, S. E. Bojesen, B. Bonanni,  
986 D. Bondavalli, A. Borg, H. Brauch, H. Brenner, I. Briceno, A. Broeks,  
987 S. Y. Brucker, T. Brüning, B. Burwinkel, S. S. Buys, H. Byers, T. Caldés,  
988 M. A. Caligo, M. Calvello, D. Campa, J. E. Castelao, J. Chang-Claude, S. J.  
989 Chanock, M. Christiaens, H. Christiansen, W. K. Chung, K. B. Claes, C. L.  
990 Clarke, S. Cornelissen, F. J. Couch, A. Cox, S. S. Cross, K. Czene, M. B.  
991 Daly, P. Devilee, O. Diez, S. M. Domchek, T. Dörk, M. Dwek, D. M. Eccles,  
992 A. B. Ekici, D. G. Evans, P. A. Fasching, J. Figueroa, L. Foretova, F. Fo-  
993 stira, E. Friedman, D. Frost, M. Gago-Dominguez, S. M. Gapstur, J. Gar-  
994 ber, J. A. García-Sáenz, M. M. Gaudet, S. A. Gayther, G. G. Giles, A. K.  
995 Godwin, M. S. Goldberg, D. E. Goldgar, A. González-Neira, M. H. Greene,  
996 J. Gronwald, P. Guénel, L. Häberle, E. Hahnen, C. A. Haiman, C. R. Hake,  
997 P. Hall, U. Hamann, E. F. Harkness, B. A. Heemskerk-Gerritsen, P. Hille-  
998 manns, F. B. Hogervorst, B. Holleczek, A. Hollestelle, M. J. Hooning, R. N.  
999 Hoover, J. L. Hopper, A. Howell, H. Huebner, P. J. Hulick, E. N. Imyani-  
1000 tov, C. Isaacs, L. Izatt, A. Jager, M. Jakimovska, A. Jakubowska, P. James,  
1001 R. Janavicius, W. Janni, E. M. John, M. E. Jones, A. Jung, R. Kaaks, P. M.  
1002 Kapoor, B. Y. Karlan, R. Keeman, S. Khan, E. Khusnutdinova, C. M. Kita-  
1003 hara, Y. D. Ko, I. Konstantopoulou, L. B. Koppert, S. Koutros, V. N. Kris-  
1004 tensen, A. V. Laenkholm, D. Lambrechts, S. C. Larsson, P. Laurent-Puig,  
1005 C. Lazaro, E. Lazarova, F. Lejbkowitz, G. Leslie, F. Lesueur, A. Lind-  
1006 blom, J. Lissowska, W. Y. Lo, J. T. Loud, J. Lubinski, A. Lukomska, R. J.  
1007 MacInnis, A. Mannermaa, M. Manoochchri, S. Manoukian, S. Margolin,  
1008 M. E. Martinez, L. Matricardi, L. McGuffog, C. McLean, N. Mebirouk,  
1009 A. Meindl, U. Menon, A. Miller, E. Mingazheva, M. Montagna, A. M. Mulli-  
1010 gan, C. Mulot, T. A. Muranen, K. L. Nathanson, S. L. Neuhausen, H. Nevan-  
1011 linna, P. Neven, W. G. Newman, F. C. Nielsen, L. Nikitina-Zake, J. Nodora,  
1012 K. Offit, E. Olah, O. I. Olopade, H. Olsson, N. Orr, L. Papi, J. Papp, T. W.  
1013 Park-Simon, M. T. Parsons, B. Peissel, A. Peixoto, B. Peshkin, P. Peterlongo,  
1014 J. Peto, K. A. Phillips, M. Piedmonte, D. Plaseska-Karanfilska, K. Prajzen-

- 1015 danc, R. Prentice, D. Prokofyeva, B. Rack, P. Radice, S. J. Ramus, J. Rantala,  
1016 M. U. Rashid, G. Rennert, H. S. Rennert, H. A. Risch, A. Romero, M. A.  
1017 Rookus, M. Rübner, T. Rüdiger, E. Saloustros, S. Sampson, D. P. Sandler,  
1018 E. J. Sawyer, M. T. Scheuner, R. K. Schmutzler, A. Schneeweiss, M. J. Schoe-  
1019 maker, B. Schöttker, P. Schürmann, L. Senter, P. Sharma, M. E. Sherman,  
1020 X. O. Shu, C. F. Singer, S. Smichkoska, P. Soucy, M. C. Southey, J. J. Spinelli,  
1021 J. Stone, D. Stoppa-Lyonnet, A. J. Swerdlow, C. I. Szabo, R. M. Tamimi,  
1022 W. J. Tapper, J. A. Taylor, M. R. Teixeira, M. B. Terry, M. Thomassen, D. L.  
1023 Thull, M. Tischkowitz, A. E. Toland, R. A. Tollenaar, I. Tomlinson, D. Torres,  
1024 M. A. Troester, T. Truong, N. Tung, M. Untch, C. M. Vachon, A. M. van den  
1025 Ouweland, L. E. van der Kolk, E. M. van Veen, E. J. vanRensburg, A. Vega,  
1026 B. Wappenschmidt, C. R. Weinberg, J. N. Weitzel, H. Wildiers, R. Winqvist,  
1027 A. Wolk, X. R. Yang, D. Yannoukakos, W. Zheng, K. K. Zorn, R. L. Milne,  
1028 P. Kraft, J. Simard, P. D. Pharoah, K. Michailidou, A. C. Antoniou, M. K.  
1029 Schmidt, G. Chenevix-Trench, D. F. Easton, N. Chatterjee, M. García-Closas,  
1030 Genome-wide association study identifies 32 novel breast cancer susceptibil-  
1031 ity loci from overall and subtype-specific analyses, *Nature Genetics* 52 (6)  
1032 (2020) 572–581. doi:10.1038/s41588-020-0609-2.  
1033 URL <https://pubmed.ncbi.nlm.nih.gov/32424353/>
- 1034 [27] L. Darrous, N. Mounier, Z. Kutalik, Simultaneous estimation of bi-  
1035 directional causal effects and heritable confounding from GWAS sum-  
1036 mary statistics, *Nature Communications* 2021 12:1 12 (1) (2021) 1–15.  
1037 doi:10.1038/s41467-021-26970-w.  
1038 URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41467-021-26970-w)  
1039 [s41467-021-26970-w](https://www.nature.com/articles/s41467-021-26970-w)
- 1040 [28] S. Burgess, G. D. Smith, N. M. Davies, F. Dudbridge, D. Gill, M. M. Gly-  
1041 mour, F. P. Hartwig, M. V. Holmes, C. Minelli, C. L. Relton, E. Theodoratou,  
1042 G. Davey Smith, N. M. Davies, F. Dudbridge, D. Gill, M. M. Glymour, F. P.  
1043 Hartwig, M. V. Holmes, C. Minelli, C. L. Relton, E. Theodoratou, Guidelines  
1044 for performing Mendelian randomization investigations, *Wellcome Open Re-*  
1045 *search* 4 (2023) 186. doi:10.12688/wellcomeopenres.15555.2.  
1046 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7384151/)  
1047 [PMC7384151/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7384151/)
- 1048 [29] V. W. Skrivankova, R. C. Richmond, B. A. Woolf, N. M. Davies, S. A. Swan-  
1049 son, T. J. Vanderweele, N. J. Timpson, J. P. Higgins, N. Dimou, C. Langen-  
1050 berg, E. W. Loder, R. M. Golub, M. Egger, G. D. Smith, J. B. Richards,  
1051 Strengthening the reporting of observational studies in epidemiology using  
1052 mendelian randomisation (STROBE-MR): explanation and elaboration, *BMJ*

- 1053 375 (10 2021). doi:10.1136/BMJ.N2233.  
1054 URL <https://www.bmj.com/content/375/bmj.n2233>
- 1055 [30] V. W. Skrivankova, R. C. Richmond, B. A. Woolf, J. Yarmolinsky, N. M.  
1056 Davies, S. A. Swanson, T. J. Vanderweele, J. P. Higgins, N. J. Timpson,  
1057 N. Dimou, C. Langenberg, R. M. Golub, E. W. Loder, V. Gallo, A. Tybjaerg-  
1058 Hansen, G. Davey Smith, M. Egger, J. B. Richards, Strengthening the  
1059 Reporting of Observational Studies in Epidemiology Using Mendelian  
1060 Randomization: The STROBE-MR Statement, *JAMA* 326 (16) (2021)  
1061 1614–1621. doi:10.1001/JAMA.2021.18236.  
1062 URL [https://jamanetwork.com/journals/jama/  
1063 fullarticle/2785494](https://jamanetwork.com/journals/jama/fullarticle/2785494)
- 1064 [31] A. F. Schmidt, C. Finan, M. Gordillo-Marañón, F. W. Asselbergs, D. F.  
1065 Freitag, R. S. Patel, B. Tyl, S. Chopade, R. Faraway, M. Zwierzyna,  
1066 A. D. Hingorani, Genetic drug target validation using Mendelian ran-  
1067 domisation, *Nature Communications* 11 (1) (2020). doi:10.1038/  
1068 s41467-020-16969-0.  
1069 URL <http://dx.doi.org/10.1038/s41467-020-16969-0>
- 1070 [32] S. Burgess, A. Butterworth, S. G. Thompson, Mendelian randomization anal-  
1071 ysis with multiple genetic variants using summarized data, *Genetic Epidemi-*  
1072 *ology* 37 (7) (2013) 658–665. doi:10.1002/gepi.21758.  
1073 URL <https://pubmed.ncbi.nlm.nih.gov/24114802/>
- 1074 [33] J. Bowden, G. Davey Smith, S. Burgess, Mendelian randomization with  
1075 invalid instruments: Effect estimation and bias detection through Egger  
1076 regression, *International Journal of Epidemiology* 44 (2) (2015) 512–525.  
1077 doi:10.1093/ije/dyv080.  
1078 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/  
1079 PMC4849733/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4849733/)
- 1080 [34] J. Bowden, G. Davey Smith, P. C. Haycock, S. Burgess, Consistent Estima-  
1081 tion in Mendelian Randomization with Some Invalid Instruments Using a  
1082 Weighted Median Estimator, *Genetic Epidemiology* 40 (4) (2016) 304–314.  
1083 doi:10.1002/gepi.21965.  
1084 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/  
1085 PMC4849733/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4849733/)
- 1086 [35] J. Bowden, F. Del Greco M, C. Minelli, Q. Zhao, D. A. Lawlor, N. A.  
1087 Sheehan, J. Thompson, G. Davey Smith, Improving the accuracy of two-  
1088 sample summary-data Mendelian randomization: Moving beyond the NOME



- 1089 assumption, *International Journal of Epidemiology* 48 (3) (2019) 728–742.  
1090 doi:10.1093/ije/dyy258.  
1091 URL [https://academic.oup.com/ije/article-abstract/  
1092 48/3/728/5251908](https://academic.oup.com/ije/article-abstract/48/3/728/5251908)
- 1093 [36] E. Sanderson, W. Spiller, J. Bowden, Testing and correcting for weak and  
1094 pleiotropic instruments in two-sample multivariable Mendelian randomiza-  
1095 tion, *Stat Med.* 40 (25) (2021) 5434–5452.  
1096 URL <https://pubmed.ncbi.nlm.nih.gov/34338327/>
- 1097 [37] D. G. Altman, P. Royston, The cost of dichotomising continuous variables,  
1098 *British Medical Journal* 332 (7549) (2006) 1080. doi:10.1136/bmj.  
1099 332.7549.1080.
- 1100 [38] B. Elsworth, T. R. Gaunt, MELODI Presto: A fast and agile tool to explore  
1101 semantic triples derived from biomedical literature, *Bioinformatics* 37 (4)  
1102 (2021) 583–585. doi:10.1093/bioinformatics/btaa726.  
1103 URL [https://academic.oup.com/bioinformatics/  
1104 advance-article/doi/10.1093/bioinformatics/btaa726/  
1105 5893950](https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/btaa726/5893950)
- 1106 [39] T. C. Rindfleisch, M. Fiszman, The interaction of domain knowledge and  
1107 linguistic structure in natural language processing: interpreting hypernymic  
1108 propositions in biomedical text, *Journal of Biomedical Informatics* 36 (6)  
1109 (2003) 462–477. doi:10.1016/J.JBI.2003.11.003.  
1110 URL <https://pubmed.ncbi.nlm.nih.gov/14759819/>
- 1111 [40] H. Kilicoglu, G. Rosemblat, M. Fiszman, D. Shin, Broad-coverage biomed-  
1112 ical relation extraction with SemRep, *BMC Bioinformatics* 21 (1) (2020)  
1113 1–28. doi:10.1186/S12859-020-3517-7/TABLES/4.  
1114 URL [https://bmcbioinformatics.biomedcentral.com/  
1115 articles/10.1186/s12859-020-3517-7](https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-020-3517-7)
- 1116 [41] N. R. Smalheiser, Rediscovering Don Swanson: The Past, Present and Future  
1117 of Literature-based Discovery, *Journal of Data and Information Science* 2 (4)  
1118 (2017) 43–64. doi:10.1515/JDIS-2017-0019.  
1119 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/  
1120 PMC5771422/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5771422/)
- 1121 [42] D. R. Swanson, Fish oil, Raynaud’s syndrome, and undiscovered public  
1122 knowledge., *Perspectives in biology and medicine* 30 (1) (1986) 7–18. doi:  
1123 10.1353/pbm.1986.0087.  
1124 URL <https://pubmed.ncbi.nlm.nih.gov/3797213/>

- 1125 [43] M. Weeber, H. Klein, L. T. De Jong-Van Den Berg, R. Vos, Using concepts in literature-based discovery: Simulating Swanson's Raynaud-fish  
1126 oil and migraine-magnesium discoveries, *Journal of the American Society for Information Science and Technology* 52 (7) (2001) 548–557.  
1127 doi:10.1002/asi.1104.  
1128 URL [https://onlinelibrary.wiley.com/doi/10.1002/](https://onlinelibrary.wiley.com/doi/10.1002/asi.1104)  
1129 [asi.1104](https://onlinelibrary.wiley.com/doi/10.1002/asi.1104)  
1130  
1131
- 1132 [44] M. Vabistsevits, G. Davey Smith, E. Sanderson, T. G. Richardson, B. Lloyd-Lewis, R. C. Richmond, Deciphering how early life adiposity influences  
1133 breast cancer risk using Mendelian randomization, *Communications Biology* 5 (1) (2022). doi:10.1038/s42003-022-03272-5.  
1134  
1135
- 1136 [45] B. Jassal, L. Matthews, G. Viteri, C. Gong, P. Lorente, A. Fabregat, K. Sidiropoulos, J. Cook, M. Gillespie, R. Haw, F. Loney, B. May, M. Milacic, K. Rothfels, C. Sevilla, V. Shamovsky, S. Shorser, T. Varusai, J. Weiser, G. Wu, L. Stein, H. Hermjakob, P. D'Eustachio, The reactome pathway  
1137 knowledgebase, *Nucleic Acids Research* 48 (D1) (2020) D498–D503. doi:  
1138 10.1093/nar/gkz1031.  
1139 URL <https://pubmed.ncbi.nlm.nih.gov/31691815/>  
1140  
1141  
1142
- 1143 [46] H. J. Baer, S. S. Tworoger, S. E. Hankinson, W. C. Willett, Body fatness at  
1144 young ages and risk of breast cancer throughout life, *American Journal of Epidemiology* 171 (11) (2010) 1183–1194. doi:10.1093/aje/kwq045.  
1145 URL [https://academic.oup.com/aje/article-abstract/](https://academic.oup.com/aje/article-abstract/171/11/1183/101419)  
1146 [171/11/1183/101419](https://academic.oup.com/aje/article-abstract/171/11/1183/101419)  
1147
- 1148 [47] A. Furer, A. Afek, A. Sommer, L. Keinan-Boker, E. Derazne, Z. Levi, D. Tzur, S. Tiosano, A. Shina, Y. Glick, J. D. Kark, A. Tirosh, G. Twig, Adolescent obesity and midlife cancer risk: a population-based cohort study  
1149 of 2·3 million adolescents in Israel, *The Lancet Diabetes and Endocrinology* 8 (3) (2020) 216–225. doi:10.1016/S2213-8587(20)30019-X.  
1150 URL [http://dx.doi.org/10.1016/S2213-8587\(20\)30019-X](http://dx.doi.org/10.1016/S2213-8587(20)30019-X)  
1151  
1152  
1153
- 1154 [48] T. G. Richardson, E. Sanderson, B. Elsworth, K. Tilling, G. Davey Smith, Use of genetic variation to separate the effects of early and later life adiposity  
1155 on disease risk: Mendelian randomisation study, *The BMJ* 369 (2020). doi:  
1156 10.1136/bmj.m1203.  
1157 URL <http://dx.doi.org/10.1136/bmj.m1203>  
1158
- 1159 [49] Y. Hao, J. Xiao, Y. Liang, X. Wu, H. Zhang, C. Xiao, L. Zhang, S. Burgess, N. Wang, X. Zhao, P. Kraft, J. Li, X. Jiang, Reassessing the causal role  
1160

- 1161 of obesity in breast cancer susceptibility: A comprehensive multivariable  
1162 Mendelian randomization investigating the distribution and timing of  
1163 exposure, *International Journal of Epidemiology* 52 (1) (2023) 58–70.  
1164 doi:10.1093/ije/dyac143.  
1165 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/  
1166 PMC7614158/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7614158/)
- 1167 [50] S. C. Dixon-Suen, S. J. Lewis, R. M. Martin, D. R. English, T. Boyle,  
1168 G. G. Giles, K. Michailidou, M. K. Bolla, Q. Wang, J. Dennis, M. Lush,  
1169 A. Investigators, T. U. Ahearn, C. B. Ambrosone, I. L. Andrulis, H. Anton-  
1170 Culver, V. Arndt, K. J. Aronson, A. Augustinsson, P. Auvinen, L. E. Freeman,  
1171 H. Becher, M. W. Beckmann, S. Behrens, M. Bermisheva, C. Blomqvist,  
1172 N. V. Bogdanova, S. E. Bojesen, B. Bonanni, H. Brenner, T. Brüning,  
1173 S. S. Buys, N. J. Camp, D. Campa, F. Canzian, J. E. Castelao, M. H.  
1174 Cessna, J. Chang-Claude, S. J. Chanock, C. L. Clarke, D. M. Conroy, F. J.  
1175 Couch, A. Cox, S. S. Cross, K. Czene, M. B. Daly, P. Devilee, T. Dörk,  
1176 M. Dwek, D. M. Eccles, A. H. Eliassen, C. Engel, M. Eriksson, D. G.  
1177 Evans, P. A. Fasching, O. Fletcher, H. Flyger, L. Fritschi, M. Gabrielson,  
1178 M. Gago-Dominguez, M. García-Closas, J. A. García-Sáenz, M. S.  
1179 Goldberg, P. Guénel, M. Gündert, E. Hahnen, C. A. Haiman, L. Häberle,  
1180 N. Håkansson, P. Hall, U. Hamann, S. N. Hart, M. Harvie, P. Hillemanns,  
1181 A. Hollestelle, M. J. Hooning, R. Hoppe, J. Hopper, A. Howell, D. J. Hunter,  
1182 A. Jakubowska, W. Janni, E. M. John, A. Jung, R. Kaaks, R. Keeman, C. M.  
1183 Kitahara, S. Koutros, P. Kraft, V. N. Kristensen, K. Kubelka-Sabit, A. W.  
1184 Kurian, J. V. Lacey, D. Lambrechts, L. Le Marchand, A. Lindblom, S. Loibl,  
1185 J. Lubiński, A. Mannermaa, M. Manoochchri, S. Margolin, M. E. Martinez,  
1186 D. Mavroudis, U. Menon, A. M. Mulligan, R. A. Murphy, N. Collabora-  
1187 tors, H. Nevanlinna, I. Nevelsteen, W. G. Newman, K. Offit, A. F. Olshan,  
1188 H. Olsson, N. Orr, A. Patel, J. Peto, D. Plaseska-Karanfilska, N. Presneau,  
1189 B. Rack, P. Radice, E. Rees-Punia, G. Rennert, H. S. Rennert, A. Romero,  
1190 E. Saloustros, D. P. Sandler, M. K. Schmidt, R. K. Schmutzler, L. Schwent-  
1191 ner, C. Scott, M. Shah, X. O. Shu, J. Simard, M. C. Southey, J. Stone,  
1192 H. Surowy, A. J. Swerdlow, R. M. Tamimi, W. J. Tapper, J. A. Taylor, M. B.  
1193 Terry, R. A. Tollenaar, M. A. Troester, T. Truong, M. Untch, C. M. Vachon,  
1194 V. Joseph, B. Wappenschmidt, C. R. Weinberg, A. Wolk, D. Yannoukakos,  
1195 W. Zheng, A. Ziogas, A. M. Dunning, P. D. Pharoah, D. F. Easton, R. L.  
1196 Milne, B. M. Lynch, Physical activity, sedentary time and breast cancer risk:  
1197 a Mendelian randomisation study, *British Journal of Sports Medicine* 56 (20)  
1198 (2022) 1157–1170. doi:10.1136/bjsports-2021-105132.  
1199 URL <https://bjsm.bmj.com/content/56/20/1157>

- 1200 [51] R. C. Richmond, E. L. Anderson, H. S. Dashti, S. E. Jones, J. M. Lane, L. B.  
1201 Strand, B. Brumpton, M. K. Rutter, A. R. Wood, K. Straif, C. L. Relton,  
1202 M. Munafò, T. M. Frayling, R. M. Martin, R. Saxena, M. N. Weedon, D. A.  
1203 Lawlor, G. D. Smith, Investigating causal relations between sleep traits and  
1204 risk of breast cancer in women: Mendelian randomisation study, *The BMJ*  
1205 365 (2019). doi:10.1136/bmj.12327.  
1206 URL <http://dx.doi.org/10.1136/bmj.12327>
- 1207 [52] C. Nowak, J. Ärnlöv, A Mendelian randomization study of the effects of  
1208 blood lipids on breast cancer risk, *Nature communications* 9 (1) (12 2018).  
1209 doi:10.1038/S41467-018-06467-9.  
1210 URL <https://pubmed.ncbi.nlm.nih.gov/30262900/>
- 1211 [53] A. Beeghly-Fadiel, N. K. Khankari, R. J. Delahanty, X. O. Shu, Y. Lu,  
1212 M. K. Schmidt, M. K. Bolla, K. Michailidou, Q. Wang, J. Dennis, D. Yan-  
1213 noukakos, A. M. Dunning, P. D. Pharoah, G. Chenevix-Trench, R. L. Milne,  
1214 D. J. Hunter, H. Per, P. Kraft, J. Simard, D. F. Easton, W. Zheng, A  
1215 Mendelian randomization analysis of circulating lipid traits and breast cancer  
1216 risk, *International Journal of Epidemiology* 49 (4) (2020) 1117–1131.  
1217 doi:10.1093/ije/dyz242.  
1218 URL <https://pubmed.ncbi.nlm.nih.gov/31872213/>
- 1219 [54] K. E. Johnson, K. M. Siewert, D. Klarin, S. M. Damrauer, K. M. Chang, P. S.  
1220 Tsao, T. L. Assimes, K. N. Maxwell, B. F. Voight, The relationship between  
1221 circulating lipids and breast cancer risk: A Mendelian randomization study,  
1222 *PLoS Medicine* 17 (9) (2020) 1–21. doi:10.1371/JOURNAL.PMED.  
1223 1003302.  
1224 URL <https://pubmed.ncbi.nlm.nih.gov/32915777/>
- 1225 [55] M. Touvier, P. Fassier, M. His, T. Norat, D. S. Chan, J. Blacher, S. Her-  
1226 cberg, P. Galan, N. Druesne-Pecollo, P. Latino-Martel, Cholesterol and breast  
1227 cancer risk: A systematic review and meta-analysis of prospective stud-  
1228 ies, *British Journal of Nutrition* 114 (3) (2015) 347–357. doi:10.1017/  
1229 S000711451500183X.  
1230 URL <https://pubmed.ncbi.nlm.nih.gov/26173770/>
- 1231 [56] L. Lindstedt, M. Lee, K. Öörni, D. Brömme, P. T. Kovanen, Cathepsins F and  
1232 S block HDL3-induced cholesterol efflux from macrophage foam cells, *Bio-  
1233 chemical and Biophysical Research Communications* 312 (4) (2003) 1019–  
1234 1024. doi:10.1016/j.bbrc.2003.11.020.  
1235 URL <https://pubmed.ncbi.nlm.nih.gov/14651973/>

- 1236 [57] N. Spielmann, D. M. Mutch, F. Rousseau, F. Tores, J. Hager, S. Bertrais,  
1237 A. Basdevant, P. Tounian, B. Dubern, P. Galan, K. Clément, Cathepsin S  
1238 genotypes are associated with Apo-A1 and HDL-cholesterol in lean and  
1239 obese French populations, *Clinical Genetics* 74 (2) (2008) 155–163. doi:  
1240 10.1111/j.1399-0004.2008.01043.x.  
1241 URL <https://pubmed.ncbi.nlm.nih.gov/18565099/>
- 1242 [58] L. Ren, J. Yi, W. Li, X. Zheng, J. Liu, J. Wang, G. Du, Apolipopro-  
1243 teins and cancer, *Cancer Medicine* 8 (16) (2019) 7032–7043.  
1244 doi:10.1002/cam4.2587.  
1245 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6853823/)  
1246 [PMC6853823/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6853823/)
- 1247 [59] O. C. Olson, J. A. Joyce, Cysteine cathepsin proteases: Regulators of cancer  
1248 progression and therapeutic response, *Nature Reviews Cancer* 15 (12) (2015)  
1249 712–729. doi:10.1038/nrc4027.  
1250 URL <https://pubmed.ncbi.nlm.nih.gov/26597527/>
- 1251 [60] M. Rudzińska, A. Parodi, S. M. Soond, A. Z. Vinarov, D. O. Korolev, A. O.  
1252 Morozov, C. Daglioglu, Y. Tutar, A. A. Zamyatin, The Role of Cysteine  
1253 Cathepsins in Cancer Progression and Drug Resistance, *International Journal*  
1254 *of Molecular Sciences* 20 (14) (7 2019). doi:10.3390/IJMS20143602.  
1255 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6678516/)  
1256 [PMC6678516/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6678516/)
- 1257 [61] K. S. Ruth, F. R. Day, J. Tyrrell, D. J. Thompson, A. R. Wood, A. Mahajan,  
1258 R. N. Beaumont, L. Wittemans, S. Martin, A. S. Busch, A. M. Erzurum-  
1259 luoglu, B. Hollis, T. A. O’Mara, M. I. McCarthy, C. Langenberg, D. F.  
1260 Easton, N. J. Wareham, S. Burgess, A. Murray, K. K. Ong, T. M. Frayling,  
1261 J. R. Perry, Using human genetics to understand the disease impacts of  
1262 testosterone in men and women, *Nature Medicine* 26 (2) (2020) 252–258.  
1263 doi:10.1038/s41591-020-0751-5.  
1264 URL [http://www.nature.com/articles/](http://www.nature.com/articles/s41591-020-0751-5)  
1265 [s41591-020-0751-5](http://www.nature.com/articles/s41591-020-0751-5)
- 1266 [62] Y. Guo, S. Warren Andersen, X. O. Shu, K. Michailidou, M. K. Bolla,  
1267 Q. Wang, M. Garcia-Closas, R. L. Milne, M. K. Schmidt, J. Chang-Claude,  
1268 A. Dunning, S. E. Bojesen, H. Ahsan, K. Aittomäki, I. L. Andrulis, H. Anton-  
1269 Culver, V. Arndt, M. W. Beckmann, A. Beeghly-Fadiel, J. Benitez, N. V.  
1270 Bogdanova, B. Bonanni, A. L. Børresen-Dale, J. Brand, H. Brauch, H. Bren-  
1271 ner, T. Brüning, B. Burwinkel, G. Casey, G. Chenevix-Trench, F. J. Couch,  
1272 A. Cox, S. S. Cross, K. Czene, P. Devilee, T. Dörk, M. Dumont, P. A.

- 1273 Fasching, J. Figueroa, D. Flesch-Janys, O. Fletcher, H. Flyger, F. Fostira,  
1274 M. Gammon, G. G. Giles, P. Guénel, C. A. Haiman, U. Hamann, M. J.  
1275 Hooning, J. L. Hopper, A. Jakubowska, F. Jasmine, M. Jenkins, E. M. John,  
1276 N. Johnson, M. E. Jones, M. Kabisch, M. Kibriya, J. A. Knight, L. B. Kop-  
1277 pert, V. M. Kosma, V. Kristensen, L. Le Marchand, E. Lee, J. Li, A. Lind-  
1278 blom, R. Luben, J. Lubinski, K. E. Malone, A. Mannermaa, S. Margolin,  
1279 F. Marme, C. McLean, H. Meijers-Heijboer, A. Meindl, S. L. Neuhausen,  
1280 H. Nevanlinna, P. Neven, J. E. Olson, J. I. Perez, B. Perkins, P. Peter-  
1281 longo, K. A. Phillips, K. Pylkäs, A. Rudolph, R. Santella, E. J. Sawyer,  
1282 R. K. Schmutzler, C. Seynaeve, M. Shah, M. J. Shrubsole, M. C. Southey,  
1283 A. J. Swerdlow, A. E. Toland, I. Tomlinson, D. Torres, T. Truong, G. Ursin,  
1284 R. B. Van Der Luijt, S. Verhoef, A. S. Whittemore, R. Winqvist, H. Zhao,  
1285 S. Zhao, P. Hall, J. Simard, P. Kraft, P. Pharoah, D. Hunter, D. F. Easton,  
1286 W. Zheng, Genetically Predicted Body Mass Index and Breast Cancer Risk:  
1287 Mendelian Randomization Analyses of Data from 145,000 Women of Euro-  
1288 pean Descent, *PLoS Medicine* 13 (8) (2016) e1002105. doi:10.1371/  
1289 journal.pmed.1002105.  
1290 URL <http://dx.plos.org/10.1371/journal.pmed.1002105>
- 1291 [63] D. Freuer, J. Linseisen, T. A. O'mara, M. Leitzmann, H. Baurecht, S. E.  
1292 Baumeister, C. Meisinger, Body Fat Distribution and Risk of Breast, En-  
1293 dometrial, and Ovarian Cancer: A Two-Sample Mendelian Randomization  
1294 Study, *Cancers* 13 (20) (10 2021). doi:10.3390/CANCERS13205053.  
1295 URL <https://pubmed.ncbi.nlm.nih.gov/34680200/>
- 1296 [64] B. L. Hayes, T. Robinson, S. Kar, K. S. Ruth, K. K. Tsilidis, T. Frayling,  
1297 A. Murray, R. M. Martin, D. A. Lawlor, R. C. Richmond, Do sex hormones  
1298 confound or mediate the effect of chronotype on breast and prostate cancer?  
1299 A Mendelian randomization study, *PLoS Genetics* 18 (1) (2022) 1–28.  
1300 doi:10.1371/journal.pgen.1009887.  
1301 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8809575)  
1302 [PMC8809575](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8809575)
- 1303 [65] M. Vithayathil, P. Carter, S. Kar, A. M. Mason, S. Burgess, S. C. Larsson,  
1304 S. Karx, A. M. Mason, S. Burgess, S. C. Larsson, Body size and compo-  
1305 sition and site-specific cancers in UK Biobank: a Mendelian randomisation  
1306 study, *PLoS Medicine* 18 (7) (2021) 1–20. doi:10.1101/2020.02.28.  
1307 970459.  
1308 URL <http://dx.doi.org/10.1371/journal.pmed.1003706>
- 1309 [66] N. Dimou, J. Yarmolinsky, E. Bouras, K. K. Tsilidis, R. M. Martin, S. J.

- 1310 Lewis, I. T. Gram, M. F. Bakker, H. Brenner, J. C. Figueiredo, R. T. Fortner,  
1311 S. B. Gruber, B. van Guelpen, L. Hsu, R. Kaaks, S. S. Kweon, Y. Lin, N. M.  
1312 Lindor, P. A. Newcomb, M. J. Sanchez, G. Severi, H. A. Tindle, R. Tu-  
1313 mino, E. Weiderpass, M. J. Gunter, N. Murphy, Causal effects of lifetime  
1314 smoking on breast and colorectal cancer risk: Mendelian randomization  
1315 study, *Cancer Epidemiology Biomarkers and Prevention* 30 (5) (2021)  
1316 953–964. doi:10.1158/1055-9965.EPI-20-1218/71015/AM/  
1317 CAUSAL-EFFECTS-OF-LIFETIME-SMOKING-ON-BREAST-AND.  
1318 URL [https://aacrjournals.org/cebpa/article/30/5/953/  
1319 670780/Causal-Effects-of-Lifetime-Smoking-on-Breast-and](https://aacrjournals.org/cebpa/article/30/5/953/670780/Causal-Effects-of-Lifetime-Smoking-on-Breast-and)
- 1320 [67] H. A. Park, S. Neumeyer, K. Michailidou, M. K. Bolla, Q. Wang, J. Dennis,  
1321 T. U. Ahearn, I. L. Andrulis, H. Anton-Culver, N. N. Antonenkova, V. Arndt,  
1322 K. J. Aronson, A. Augustinsson, A. Baten, L. E. Beane Freeman, H. Becher,  
1323 M. W. Beckmann, S. Behrens, J. Benitez, M. Bermisheva, N. V. Bogdanova,  
1324 S. E. Bojesen, H. Brauch, H. Brenner, S. Y. Brucker, B. Burwinkel, D. Campa,  
1325 F. Canzian, J. E. Castelao, S. J. Chanock, G. Chenevix-Trench, C. L. Clarke,  
1326 A. L. Børresen-Dale, G. I. Grenaker Alnæs, K. K. Sahlberg, L. Ottestad,  
1327 R. Kåresen, E. Schlichting, M. M. Holmen, T. Sauer, V. Haakensen, O. En-  
1328 gebråten, B. Naume, A. Fosså, C. E. Kiserud, K. V. Reinertsen, A. Helland,  
1329 M. Riis, J. Geisler, D. M. Conroy, F. J. Couch, A. Cox, S. S. Cross, K. Czene,  
1330 M. B. Daly, P. Devilee, T. Dörk, I. dos Santos-Silva, M. Dwek, D. M. Eccles,  
1331 A. H. Eliassen, C. Engel, M. Eriksson, D. G. Evans, P. A. Fasching, H. Flyger,  
1332 L. Fritschi, M. García-Closas, J. A. García-Sáenz, M. M. Gaudet, G. G. Giles,  
1333 G. Glendon, M. S. Goldberg, D. E. Goldgar, A. González-Neira, M. Grip,  
1334 P. Guénel, E. Hahnen, C. A. Haiman, N. Håkansson, P. Hall, U. Hamann,  
1335 S. Han, E. F. Harkness, S. N. Hart, W. He, B. A. Heemskerk-Gerritsen, J. L.  
1336 Hopper, D. J. Hunter, C. Clarke, D. Marsh, R. Scott, R. Baxter, D. Yip,  
1337 J. Carpenter, A. Davis, N. Pathmanathan, P. Simpson, D. Graham, M. Sach-  
1338 chithananthan, D. Amor, L. Andrews, Y. Antill, R. Balleine, J. Beesley,  
1339 I. Bennett, M. Bogwitz, L. Botes, M. Brennan, M. Brown, M. Buckley,  
1340 J. Burke, P. Butow, L. Caldon, I. Campbell, D. Chauhan, M. Chauhan,  
1341 A. Christian, P. Cohen, A. Colley, A. Crook, J. Cui, M. Cummings, S. J.  
1342 Dawson, A. DeFazio, M. Delatycki, R. Dickson, J. Dixon, T. Edkins, S. Ed-  
1343 wards, G. Farshid, A. Fellows, G. Fenton, M. Field, J. Flanagan, P. Fong,  
1344 L. Forrest, S. Fox, J. French, M. Friedlander, C. Gaff, M. Gattas, P. George,  
1345 S. Greening, M. Harris, S. Hart, N. Hayward, C. Hoskins, C. Hunt, P. James,  
1346 M. Jenkins, A. Kidd, J. Kirk, J. Koehler, J. Kollias, S. Lakhani, M. Lawrence,  
1347 G. Lindeman, L. Lipton, L. Lobb, G. Mann, D. Marsh, S. A. McLach-  
1348 lan, B. Meiser, S. Nightingale, S. O’Connell, S. O’Sullivan, D. G. Or-

- 1349 tega, N. Pachter, B. Patterson, A. Pearn, K. Phillips, E. Pieper, E. Rickard,  
1350 B. Robinson, M. Saleh, E. Salisbury, C. Saunders, J. Saunus, R. Scott,  
1351 C. Scott, A. Sexton, A. Shelling, P. Simpson, A. Spurdle, J. Taylor, R. Taylor,  
1352 H. Thorne, A. Trainer, K. Tucker, J. Visvader, L. Walker, R. Williams, I. Win-  
1353 ship, M. A. Young, A. Jager, A. Jakubowska, E. M. John, A. Jung, R. Kaaks,  
1354 P. M. Kapoor, R. Keeman, E. Khusnutdinova, C. M. Kitahara, L. B. Kop-  
1355 pert, S. Koutros, V. N. Kristensen, A. W. Kurian, J. Lacey, D. Lambrechts,  
1356 L. Le Marchand, W. Y. Lo, J. Lubiński, A. Mannermaa, M. Manoochchri,  
1357 S. Margolin, M. E. Martinez, D. Mavroudis, A. Meindl, U. Menon, R. L.  
1358 Milne, T. A. Muranen, H. Nevanlinna, W. G. Newman, B. G. Nordestgaard,  
1359 K. Offit, A. F. Olshan, T. W. Park-Simon, P. Peterlongo, J. Peto, D. Plaseska-  
1360 Karanfiliska, P. Radice, G. Rennert, H. S. Rennert, A. Romero, E. Salous-  
1361 tros, E. J. Sawyer, M. K. Schmidt, R. K. Schmutzler, M. J. Schoemaker,  
1362 L. Schwentner, C. Scott, M. Shah, X. O. Shu, J. Simard, A. Smeets, M. C.  
1363 Southey, J. J. Spinelli, V. Stevens, A. J. Swerdlow, R. M. Tamimi, W. J. Tap-  
1364 per, J. A. Taylor, M. B. Terry, I. Tomlinson, M. A. Troester, T. Truong, C. M.  
1365 Vachon, E. M. van Veen, J. Vijai, S. Wang, C. Wendt, R. Winqvist, A. Wolk,  
1366 A. Ziogas, A. M. Dunning, P. D. Pharoah, D. F. Easton, W. Zheng, P. Kraft,  
1367 J. Chang-Claude, Mendelian randomisation study of smoking exposure in  
1368 relation to breast cancer risk, *British journal of cancer* 125 (8) (2021) 1135–  
1369 1145. doi:10.1038/S41416-021-01432-8.  
1370 URL <https://pubmed.ncbi.nlm.nih.gov/34341517/>
- 1371 [68] S. C. Larsson, A. M. Mason, S. Kar, M. Vithayathil, P. Carter, J. A. Baron,  
1372 K. Michaëlsson, S. Burgess, Genetically proxied milk consumption and risk  
1373 of colorectal, bladder, breast, and prostate cancer: a two-sample Mendelian  
1374 randomization study, *BMC Medicine* 18 (1) (2020) 1–7. doi:10.1186/  
1375 S12916-020-01839-9/FIGURES/1.  
1376 URL [https://bmcmmedicine.biomedcentral.com/articles/  
1377 10.1186/s12916-020-01839-9](https://bmcmmedicine.biomedcentral.com/articles/10.1186/s12916-020-01839-9)
- 1378 [69] A. L. Lumsden, A. Mulugeta, E. Hyppönen, Milk consump-  
1379 tion and risk of twelve cancers: A large-scale observational and  
1380 Mendelian randomisation study, *Clinical Nutrition* 42 (1) (2023) 1–8.  
1381 doi:10.1016/j.clnu.2022.11.006.  
1382 URL [http://www.clinicalnutritionjournal.com/  
1383 article/S0261561422003910/fulltext](http://www.clinicalnutritionjournal.com/article/S0261561422003910/fulltext)
- 1384 [70] N. Papadimitriou, N. Dimou, K. K. Tsilidis, B. Banbury, R. M. Martin, S. J.  
1385 Lewis, N. Kazmi, T. M. Robinson, D. Albanes, K. Aleksandrova, S. I. Berndt,  
1386 D. Timothy Bishop, H. Brenner, D. D. Buchanan, B. Bueno-de Mesquita, P. T.



- 1387 Campbell, S. Castellví-Bel, A. T. Chan, J. Chang-Claude, M. Ellingjord-Dale,  
1388 J. C. Figueiredo, S. J. Gallinger, G. G. Giles, E. Giovannucci, S. B. Gruber,  
1389 A. Gsur, J. Hampe, H. Hampel, S. Harlid, T. A. Harrison, M. Hoffmeister,  
1390 J. L. Hopper, L. Hsu, J. María Huerta, J. R. Huyghe, M. A. Jenkins, T. O.  
1391 Keku, T. Kühn, C. La Vecchia, L. Le Marchand, C. I. Li, L. Li, A. Lind-  
1392 blom, N. M. Lindor, B. Lynch, S. D. Markowitz, G. Masala, A. M. May,  
1393 R. Milne, E. Monninkhof, L. Moreno, V. Moreno, P. A. Newcomb, K. Of-  
1394 fit, V. Perduca, P. D. Pharoah, E. A. Platz, J. D. Potter, G. Rennert, E. Ri-  
1395 boli, M. J. Sánchez, S. L. Schmit, R. E. Schoen, G. Severi, S. Sieri, M. L.  
1396 Slattery, M. Song, C. M. Tangen, S. N. Thibodeau, R. C. Travis, A. Tri-  
1397 chopoulou, C. M. Ulrich, F. J. van Duijnhoven, B. Van Guelpen, P. Vod-  
1398 icka, E. White, A. Wolk, M. O. Woods, A. H. Wu, U. Peters, M. J. Gunter,  
1399 N. Murphy, Physical activity and risks of breast and colorectal cancer: a  
1400 Mendelian randomisation analysis, *Nature communications* 11 (1) (2020).  
1401 doi:10.1038/S41467-020-14389-8.  
1402 URL <https://pubmed.ncbi.nlm.nih.gov/32001714/>
- 1403 [71] H. Li, L. Hou, Y. Yu, X. Sun, X. Liu, Y. Yu, S. Wu, Y. He, Y. Wu, L. He,  
1404 F. Xue, Lipids, Anthropometric Measures, Smoking and Physical Activity  
1405 Mediate the Causal Pathway From Education to Breast Cancer in Women:  
1406 A Mendelian Randomization Study, *Journal of Breast Cancer* 24 (6) (2021)  
1407 504–519. doi:10.4048/jbc.2021.24.e53.  
1408 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/  
1409 PMC8724372/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8724372/)
- 1410 [72] Y. Feng, R. Wang, C. Li, X. Cai, Z. Huo, Z. Liu, F. Ge, C. Huang, Y. Lu,  
1411 R. Zhong, J. Li, B. Cheng, H. Liang, S. Xiong, X. Mao, Y. Chen, R. Lan,  
1412 Y. Wen, H. Peng, Y. Jiang, Z. Su, X. Wu, J. He, W. Liang, Causal ef-  
1413 fects of genetically determined metabolites on cancers included lung, breast,  
1414 ovarian cancer, and glioma: a Mendelian randomization study, *Transla-  
1415 tional lung cancer research* 11 (7) (2022) 1302–1314. doi:10.21037/  
1416 TLCR-22-34.  
1417 URL <https://pubmed.ncbi.nlm.nih.gov/35958335/>
- 1418 [73] T. Robinson, R. M. Martin, J. Yarmolinsky, Mendelian randomisation anal-  
1419 ysis of circulating adipokines and C-reactive protein on breast cancer risk,  
1420 *International journal of cancer* 147 (6) (2020) 1597–1603. doi:10.1002/  
1421 IJC.32947.  
1422 URL <https://pubmed.ncbi.nlm.nih.gov/32134113/>
- 1423 [74] E. Bouras, V. Karhunen, D. Gill, J. Huang, P. C. Haycock, M. J. Gunter,

- 1424 M. Johansson, P. Brennan, T. Key, S. J. Lewis, R. M. Martin, N. Mur-  
1425 phy, E. A. Platz, R. Travis, J. Yarmolinsky, V. Zuber, P. Martin, M. Kat-  
1426 soulis, H. Freisling, T. H. Nøst, M. B. Schulze, L. Dossus, R. J. Hung, C. I.  
1427 Amos, A. Ahola-Olli, S. Palaniswamy, M. Männikkö, J. Auvinen, K. H.  
1428 Herzig, S. Keinänen-Kiukaanniemi, T. Lehtimäki, V. Salomaa, O. Raitakari,  
1429 M. Salmi, S. Jalkanen, CRUK, CAPS, PEGASUS, M. R. Jarvelin, A. De-  
1430 hghan, K. K. Tsilidis, Circulating inflammatory cytokines and risk of five  
1431 cancers: a Mendelian randomization analysis, *BMC Medicine* 20 (1) (2022)  
1432 1–15. doi:10.1186/s12916-021-02193-0.  
1433 URL [https://bmcmmedicine.biomedcentral.com/articles/](https://bmcmmedicine.biomedcentral.com/articles/10.1186/s12916-021-02193-0)  
1434 [10.1186/s12916-021-02193-0](https://bmcmmedicine.biomedcentral.com/articles/10.1186/s12916-021-02193-0)
- 1435 [75] J. Yarmolinsky, J. W. Robinson, D. Mariosa, V. Karhunen, J. Huang, N. Di-  
1436 mou, N. Murphy, K. Burrows, E. Bouras, K. Smith-Byrne, S. J. Lewis, T. E.  
1437 Galesloot, L. A. Kiemeny, S. Vermeulen, P. Martin, D. Albanes, L. Hou,  
1438 P. A. Newcomb, E. White, A. Wolk, A. H. Wu, L. Le Marchand, A. I. Phipps,  
1439 D. D. Buchanan, M. T. Landi, V. Stevens, Y. Wang, D. Albanes, N. Caporaso,  
1440 P. Brennan, C. I. Amos, S. Shete, R. J. Hung, H. Bickeböller, A. Risch,  
1441 R. Houlston, S. Lam, A. Tardon, C. Chen, S. E. Bojesen, M. Johansson, H. E.  
1442 Wichmann, D. Christiani, G. Rennert, S. Arnold, J. K. Field, L. Le Marchand,  
1443 O. Melander, H. Brunnström, G. Liu, A. Andrew, H. Shen, S. Zienolddiny,  
1444 K. Grankvist, M. Johansson, M. D. Teare, Y. C. Hong, J. M. Yuan, P. Lazarus,  
1445 M. B. Schabath, M. C. Aldrich, R. A. Eeles, C. A. Haiman, Z. Kote-Jarai,  
1446 F. R. Schumacher, S. Benlloch, A. A. Al Olama, K. R. Muir, S. I. Berndt,  
1447 D. V. Conti, F. Wiklund, S. Chanock, C. M. Tangen, J. Batra, J. A. Clements,  
1448 H. Grönberg, N. Pashayan, J. Schleutker, S. J. Weinstein, C. M. West, L. A.  
1449 Mucci, G. Cancel-Tassin, S. Koutros, K. D. Sørensen, E. M. Grindedal,  
1450 D. E. Neal, F. C. Hamdy, J. L. Donovan, R. C. Travis, R. J. Hamilton, S. A.  
1451 Ingles, B. S. Rosenstein, Y. J. Lu, G. G. Giles, R. J. MacInnis, A. S. Kibel,  
1452 A. Vega, M. Kogevinas, K. L. Penney, J. Y. Park, J. L. Stanfrod, C. Cybulski,  
1453 B. G. Nordestgaard, S. F. Nielsen, H. Brenner, C. Maier, C. J. Logothetis,  
1454 E. M. John, M. R. Teixeira, S. L. Neuhausen, K. De Ruyck, A. Razack, L. F.  
1455 Newcomb, D. Lessel, R. Kaneva, N. Usmani, F. Claessens, P. A. Townsend,  
1456 J. E. Castelao, M. J. Roobol, F. Menegaux, K. T. Khaw, L. Cannon-Albright,  
1457 H. Pandha, S. N. Thibodeau, D. J. Hunter, P. Kraft, W. J. Blot, E. Riboli,  
1458 S. S. Zhao, D. Gill, S. J. Chanock, M. P. Purdue, G. Davey Smith, K. H.  
1459 Herzig, M. R. Jarvelin, C. I. Amos, A. Dehghan, M. J. Gunter, K. K. Tsilidis,  
1460 R. M. Martin, Association between circulating inflammatory markers and  
1461 adult cancer risk: a Mendelian randomization analysis, *eBioMedicine* 100  
1462 (2 2024). doi:10.1016/J.EBIOM.2024.104991/ATTACHMENT/

- 1463 0C4432FD-1E6E-480A-A8CA-D4C73CF8C81C/MMC2.DOCX.  
1464 URL [http://www.thelancet.com/article/  
1465 S2352396424000264/fulltext](http://www.thelancet.com/article/S2352396424000264/fulltext)
- 1466 [76] Z. Zhao, Q. Cao, M. Zhu, C. Wang, X. Lu, Causal relationships between  
1467 serum matrix metalloproteinases and estrogen receptor-negative breast  
1468 cancer: a bidirectional mendelian randomization study, *Scientific Reports*  
1469 2023 13:1 13 (1) (2023) 1–25. doi:10.1038/s41598-023-34200-0.  
1470 URL [https://www.nature.com/articles/  
1471 s41598-023-34200-0](https://www.nature.com/articles/s41598-023-34200-0)
- 1472 [77] X. Shu, L. Wu, N. K. Khankari, X. O. Shu, T. J. Wang, K. Michailidou,  
1473 M. K. Bolla, Q. Wang, J. Dennis, R. L. Milne, M. K. Schmidt, P. D. Pharoah,  
1474 I. L. Andrulis, D. J. Hunter, J. Simard, D. F. Easton, W. Zheng, Associations  
1475 of obesity and circulating insulin and glucose with breast cancer risk: A  
1476 Mendelian randomization analysis, *International Journal of Epidemiology*  
1477 48 (3) (2019) 795–806. doi:10.1093/ije/dyy201.  
1478 URL [http://ccge.medschl.cam.ac.uk/research/  
1479 consortia/](http://ccge.medschl.cam.ac.uk/research/consortia/)
- 1480 [78] K. Smith-Byrne, A. Hedman, M. Dimitriou, T. Desai, A. V. Sokolov, H. B.  
1481 Schioth, M. Koprulu, M. Pietzner, C. Langenberg, J. Atkins, R. Cortez,  
1482 J. McKay, P. Brennan, S. Zhou, B. J. Richards, J. Yarmolinsky, R. M.  
1483 Martin, J. Borlido, X. J. Mu, A. Butterworth, X. Shen, J. Wilson, T. L.  
1484 Assimes, R. J. Hung, C. Amos, M. Purdue, N. Rothman, S. Chanock,  
1485 R. C. Travis, M. Johansson, A. Målarstig, Identifying therapeutic targets for  
1486 cancer: 2,094 circulating proteins and risk of nine cancers, *medRxiv* (2023)  
1487 2023.05.05.23289547doi:10.1101/2023.05.05.23289547.  
1488 URL [https://www.medrxiv.org/content/10.1101/2023.  
1489 05.05.23289547v1](https://www.medrxiv.org/content/10.1101/2023.05.05.23289547v1)
- 1490 [79] F. Chen, W. Wen, J. Long, X. Shu, Y. Yang, X. Shu, W. Zheng, Mendelian  
1491 randomization analyses of 23 known and suspected risk factors and biomark-  
1492 ers for breast cancer overall and by molecular subtypes, *International Journal*  
1493 *of Cancer* (4 2022). doi:10.1002/IJC.34026.  
1494 URL [https://onlinelibrary.wiley.com/doi/10.1002/  
1495 ijc.34026](https://onlinelibrary.wiley.com/doi/10.1002/ijc.34026)
- 1496 [80] M. Escala-Garcia, A. Morra, S. Canisius, J. Chang-Claude, S. Kar, W. Zheng,  
1497 S. E. Bojesen, D. Easton, P. D. P. Pharoah, M. K. Schmidt, Breast cancer risk  
1498 factors and their effects on survival: a Mendelian randomisation study, *BMC*

- 1499 Medicine 18 (1) (2020) 1–10. doi:10.1186/s12916-020-01797-2.  
1500 URL <https://doi.org/10.1186/s12916-020-01797-2>
- 1501 [81] D. A. Lawlor, Commentary: Two-sample Mendelian randomization: Oppor-  
1502 tunities and challenges, *International Journal of Epidemiology* 45 (3) (2016)  
1503 908–915. doi:10.1093/ije/dyw127.  
1504 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/  
1505 PMC5005949/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5005949/)
- 1506 [82] M. R. Munafo, J. Brown, M. Hefler, G. Davey Smith, Man-  
1507 aging the exponential growth of mendelian randomization stud-  
1508 ies, *Tobacco Control* 33 (5) (2024) 559–560. arXiv:<https://tobaccocontrol.bmj.com/content/33/5/559.full.pdf>,  
1509 doi:10.1136/tc-2024-058987.  
1510 URL <https://tobaccocontrol.bmj.com/content/33/5/559>
- 1512 [83] G. D. Smith, S. Ebrahim, Mendelian randomisation at 20 years: how can it  
1513 avoid hubris, while achieving more?, *The Lancet Diabetes and Endocrinology*  
1514 8587 (23) (2023) 10–12. doi:10.1016/S2213-8587(23)00348-0.  
1515 URL [https://doi.org/10.1016/S2213-8587\(23\)00348-0](https://doi.org/10.1016/S2213-8587(23)00348-0)
- 1516 [84] S. Stender, H. Gellert-Kristensen, G. D. Smith, Reclaiming mendelian  
1517 randomization from the deluge of papers and misleading findings,  
1518 *Lipids in Health and Disease* 23 (1) (2024) 286. doi:10.1186/  
1519 s12944-024-02284-w.  
1520 URL <https://doi.org/10.1186/s12944-024-02284-w>
- 1521 [85] A. Buniello, J. A. MacArthur, M. Cerezo, L. W. Harris, J. Hayhurst,  
1522 C. Malangone, A. McMahon, J. Morales, E. Mountjoy, E. Sollis, D. Su-  
1523 veiges, O. Vrousitou, P. L. Whetzel, R. Amode, J. A. Guillen, H. S. Riat,  
1524 S. J. Trevanion, P. Hall, H. Junkins, P. Flicek, T. Burdett, L. A. Hindorf,  
1525 F. Cunningham, H. Parkinson, The NHGRI-EBI GWAS Catalog of pub-  
1526 lished genome-wide association studies, targeted arrays and summary statis-  
1527 tics 2019, *Nucleic Acids Research* 47 (D1) (2019) D1005–D1012. doi:  
1528 10.1093/nar/gky1120.  
1529 URL <https://pubmed.ncbi.nlm.nih.gov/30445434/>
- 1530 [86] B. B. Sun, J. Chiou, M. Traylor, C. Benner, Y. H. Hsu, T. G. Richardson,  
1531 P. Surendran, A. Mahajan, C. Robins, S. G. Vasequez-Grinnell, L. Hou, E. M.  
1532 Kvikstad, O. S. Burren, J. Davitt, K. L. Ferber, C. E. Gillies, A. K. Hed-  
1533 man, S. Hu, T. Lin, R. Mikkilineni, R. K. Pendergrass, C. Pickering, B. Prins,  
1534 D. Baird, C. Y. Chen, L. D. Ward, A. M. Deaton, S. Welsh, C. M. Willis,

- 1535 N. Lehner, M. Arnold, M. A. Wörheide, K. Suhre, G. Kastenmüller, A. Sethi,  
1536 M. Cule, A. Raj, H. M. Kang, L. Burkitt-Gray, E. Melamud, M. H. Black,  
1537 E. B. Fauman, J. M. Howson, H. M. Kang, M. I. McCarthy, P. Nioi, S. Petrovski,  
1538 R. A. Scott, E. N. Smith, S. Szalma, D. M. Waterworth, L. J. Mitnau,  
1539 J. D. Szustakowski, B. W. Gibson, M. R. Miller, C. D. Whelan, Plasma  
1540 proteomic associations with genetics and health in the UK Biobank, *Nature*  
1541 622 (7982) (2023) 329–338. doi:10.1038/s41586-023-06592-6.  
1542 URL <https://pubmed.ncbi.nlm.nih.gov/37794186/>
- 1543 [87] H. Julkunen, A. Cichońska, M. Tiainen, H. Koskela, K. Nybo, V. Mäkelä,  
1544 J. Nokso-Koivisto, K. Kristiansson, M. Perola, V. Salomaa, P. Jousilahti,  
1545 A. Lundqvist, A. J. Kangas, P. Soininen, J. C. Barrett, P. Würtz, Atlas  
1546 of plasma NMR biomarkers for health and disease in 118,461 individ-  
1547 uals from the UK Biobank, *Nature Communications* 14 (1) (12 2023).  
1548 doi:10.1038/s41467-023-36231-7.  
1549 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9898515/)  
1550 [PMC9898515/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9898515/)
- 1551 [88] E. Moreau, Literature-based discovery: Addressing the issue of the  
1552 subpar evaluation methodology, *Bioinformatics* 39 (2) (2023) 1–3.  
1553 doi:10.1093/bioinformatics/btad090.  
1554 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9945845/)  
1555 [PMC9945845/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9945845/)
- 1556 [89] E. Moreau, O. Hardiman, M. Heverin, D. O’Sullivan, Mining impactful dis-  
1557 coveries from the biomedical literature, *bioRxiv* (2022) 2022.10.28.514184.  
1558 URL [https://www.biorxiv.org/content/10.1101/2022.](https://www.biorxiv.org/content/10.1101/2022.10.28.514184v1)  
1559 [10.28.514184v1](https://www.biorxiv.org/content/10.1101/2022.10.28.514184v1)
- 1560 [90] O. Bodenreider, The Unified Medical Language System (UMLS): Integrating  
1561 biomedical terminology, *Nucleic Acids Research* 32 (DATABASE ISS.)  
1562 (2004) 267–270. doi:10.1093/nar/gkh061.  
1563 URL [https://www.ncbi.nlm.nih.gov/pmc/articles/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC308795/)  
1564 [PMC308795/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC308795/)
- 1565 [91] K. H. Wade, J. Yarmolinsky, E. Giovannucci, S. J. Lewis, I. Y. Millwood,  
1566 M. R. Munafò, F. Meddens, K. Burrows, J. A. Bell, N. M. Davies, D. Mar-  
1567 iosa, N. Kanerva, E. E. Vincent, K. Smith-Byrne, F. Guida, M. J. Gunter,  
1568 E. Sanderson, F. Dudbridge, S. Burgess, M. C. Cornelis, T. G. Richardson,  
1569 M. C. Borges, J. Bowden, G. Hemani, Y. Cho, W. Spiller, R. C. Richmond,  
1570 A. R. Carter, R. Langdon, D. A. Lawlor, R. G. Walters, K. S. Vimalaswaran,  
1571 A. Anderson, M. R. Sandu, K. Tilling, G. Davey Smith, R. M. Martin, C. L.

- 1572 Relton, Applying Mendelian randomization to appraise causality in relation-  
1573 ships between nutrition and cancer, *Cancer Causes & Control* 33 (2022).  
1574 doi:10.1007/s10552-022-01562-1.  
1575 URL <https://doi.org/10.1007/s10552-022-01562-1>
- 1576 [92] P. C. Haycock, S. Burgess, K. H. Wade, J. Bowden, C. Relton, G. D. Smith,  
1577 Best (but oft-forgotten) practices: the design, analysis, and interpretation of  
1578 Mendelian randomization studies (4 2016). doi:10.3945/ajcn.115.  
1579 118216.  
1580 URL <https://pubmed.ncbi.nlm.nih.gov/26961927/>
- 1581 [93] M. C. Cornelis, T. Kacprowski, C. Menni, S. Gustafsson, E. Pivin,  
1582 J. Adamski, A. Artati, C. B. Eap, G. Ehret, N. Friedrich, A. Ganna, I. Gues-  
1583 sous, G. Homuth, L. Lind, P. K. Magnusson, M. Mangino, N. L. Pedersen,  
1584 M. Pietzner, K. Suhre, H. Völzke, M. Bochud, T. D. Spector, H. J. Grabe,  
1585 E. Ingelsson, M. Burnier, O. Devuyst, P. Y. Martin, M. Mohaupt, F. Paccaud,  
1586 A. Pechere-Bertschi, B. Vogt, D. Ackermann, B. Ponte, M. Pruijm, Genome-  
1587 wide association study of caffeine metabolites provides new insights to caf-  
1588 feine metabolism and dietary caffeine-consumption behavior, *Human molec-  
1589 ular genetics* 25 (24) (2016) 5472–5482. doi:10.1093/HMG/DDW334.  
1590 URL <https://pubmed.ncbi.nlm.nih.gov/27702941/>
- 1591 [94] S. E. Steck, E. A. Murphy, Dietary patterns and cancer risk, *Nature Reviews*  
1592 *Cancer* 20 (2) (2020) 125–138. doi:10.1038/s41568-019-0227-4.  
1593 URL <http://dx.doi.org/10.1038/s41568-019-0227-4>
- 1594 [95] N. I. Pirastu, C. McDonnell, E. J. Grzeszkowiak ID, N. I. Mounier, F. Imamu-  
1595 raID, J. MerinoID, F. R. DayID, J. ZhengID, N. TabaID, M. Pina ConcasID,  
1596 L. RepettoID, K. A. KentistouID, A. RobinoID, T. EskoID, P. K. JoshiID,  
1597 K. FischerID, K. K. OngID, T. R. GauntID, Z. Kutalik, J. R. B Perry, J. F.  
1598 WilsonID, Using genetic variation to disentangle the complex relationship  
1599 between food intake and health outcomes, *PLOS Genetics* 18 (6) (2022)  
1600 e1010162. doi:10.1371/JOURNAL.PGEN.1010162.  
1601 URL [https://journals.plos.org/plosgenetics/article?  
1602 id=10.1371/journal.pgen.1010162](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1010162)
- 1603 [96] C. Jin, R. Li, T. Deng, Z. Lin, H. Li, Y. Yang, Q. Su, J. Wang, Y. Yang,  
1604 J. Wang, G. Chen, Y. Wang, Association between dried fruit intake and  
1605 pan-cancers incidence risk: A two-sample Mendelian randomization study,  
1606 *Frontiers in Nutrition* 9 (July) (2022) 1–14. doi:10.3389/fnut.2022.  
1607 899137.  
1608 URL <https://pubmed.ncbi.nlm.nih.gov/30696460/>

- 1609 [97] H. K. Neilson, M. S. Farris, C. R. Stone, M. M. Vaska, D. R. Brenner, C. M.  
1610 Friedenreich, Moderate-vigorous recreational physical activity and breast  
1611 cancer risk, stratified by menopause status: A systematic review and meta-  
1612 analysis (10 2016). doi:10.1097/GME.0000000000000745.  
1613 URL <https://pubmed.ncbi.nlm.nih.gov/27779567/>
- 1614 [98] F. Chong, Y. Wang, M. Song, Q. Sun, W. Xie, C. Song, Sedentary behavior  
1615 and risk of breast cancer: a dose-response meta-analysis from prospective  
1616 studies, *Breast cancer (Tokyo, Japan)* 28 (1) (2021) 48–59. doi:10.1007/  
1617 S12282-020-01126-8.  
1618 URL <https://pubmed.ncbi.nlm.nih.gov/32607943/>