

# 1 Simultaneous estimation of bi-directional causal effects and 2 heritable confounding from GWAS summary statistics

3 Liza Darrous<sup>1,2</sup>\*, Ninon Mounier<sup>1,2,\*</sup>, Zoltán Kutalik<sup>1,2,3</sup>†

## 4 Abstract

5 Mendelian Randomisation (MR), an increasingly popular method that estimates the  
6 causal effects of risk factors on complex human traits, has seen several extensions that relax  
7 its basic assumptions. However, most of these extensions suffer from two major limitations;  
8 their under-exploitation of genome-wide markers, and sensitivity to the presence of a heri-  
9 table confounder of the exposure-outcome relationship. To overcome these limitations, we  
10 propose a Latent Heritable Confounder MR (LHC-MR) method applicable to association  
11 summary statistics, which estimates bi-directional causal effects, direct heritabilities, and  
12 confounder effects while accounting for sample overlap. We demonstrate that LHC-MR out-  
13 performs several existing MR methods in a wide range of simulation settings and apply it to  
14 summary statistics of 13 complex traits. Besides several concordant results, LHC-MR un-  
15 ravelled new mechanisms (how being diagnosed for certain diseases might lead to improved  
16 lifestyle) and revealed new causal effects (e.g. HDL cholesterol being protective against high  
17 systolic blood pressure), hidden from standard MR methods due to a heritable confounder of  
18 opposite direction. Phenome-wide MR search suggested that the confounders indicated by  
19 LHC-MR for the birth weight-diabetes pair are likely to be obesity traits. Finally, LHC-MR  
20 results indicated that genetic correlations are predominantly driven by bi-directional causal  
21 effects and much less so by heritable confounders.

---

<sup>1</sup>University Center for Primary Care and Public Health, University of Lausanne, 1010, Switzerland

<sup>2</sup>Swiss Institute of Bioinformatics, Lausanne, 1015, Switzerland

<sup>3</sup>Genetics of Complex Traits, University of Exeter Medical School, University of Exeter, UK

\*These authors contributed equally to this work.

†Correspondence should be addressed to [zoltan.kutalik@unil.ch](mailto:zoltan.kutalik@unil.ch)

**NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

## 22 1 Introduction

23 The identification of frequent risk factors and the quantification of their impact on common  
24 diseases is a central quest for public health policy makers. Epidemiological studies aim to address  
25 this issue, but they are most often based on observational data due to its abundance over the  
26 years. Despite major methodological advances, a large majority of such studies have inherent  
27 limitations and suffer from confounding and reverse causation<sup>1,2</sup>. For these reasons, many of  
28 the reported associations found in classical epidemiological studies are mere correlates of disease  
29 risk, rather than causal factors directly involved in disease progression. Due to this limitation,  
30 additional evidence is required before developing public health interventions in a bid to reduce  
31 the future burden of diseases. While well-designed and carefully conducted randomised control  
32 trials (RCTs) remain the gold standard for causal inference, they are exceedingly expensive,  
33 time-consuming, may not be feasible for ethical reasons, and have high failure rates<sup>3,4</sup>.

34 Mendelian randomisation (MR), a natural genetic counterpart to RCTs, is an instrumental  
35 variable (IV) technique used to infer the strength of a causal relationship between a risk factor  
36 ( $X$ ) and an outcome ( $Y$ )<sup>5</sup>. To do so, it uses genetic variants ( $G$ ) as instruments and relies  
37 on three major assumptions (see Figure S1): (1) Relevance –  $G$  is robustly associated with  
38 the exposure. (2) Exchangeability –  $G$  is not associated with any confounder of the exposure-  
39 outcome relationship. (3) Exclusion restriction –  $G$  is independent of the outcome conditional  
40 on the exposure and all confounders of the exposure-outcome relationship (i.e. the only path  
41 between the instrument and the outcome is via the exposure).

42 The advantage of the MR approach is that for most heritable exposures, dozens (if not hundreds)  
43 of genetic instruments are known to date thanks to well-powered genome-wide association studies  
44 (GWASs). Each instrument can provide a causal effect estimate, which can be combined with  
45 others, by using an inverse variance-weighting (IVW) scheme (e.g. Burgess *et al.*<sup>6</sup>). However,  
46 the last assumption is particularly problematic, because genetic variants tend to be pleiotropic,  
47 i.e. exert effect on multiple traits independently. Still, it can be shown that if the instrument  
48 strength is independent of the direct effect on the outcome (InSIDE assumption) and the direct  
49 effects are on average zero, IVW-based methods will still yield consistent estimates. Methods,  
50 such as MR-Egger<sup>7</sup>, produce consistent estimates even if direct effects are allowed to have a  
51 non-zero offset. The third assumption can be further reduced to assuming that  $> 50\%$  of the  
52 instruments (or in terms of their weight) are valid (median-based estimators<sup>8</sup>) or that zero-  
53 pleiotropy instruments are the most frequent (mode-based estimators<sup>9</sup>).

54 The InSIDE assumption (i.e. horizontal pleiotropic effects ( $G \rightarrow Y$ ) are independent of the  
55 direct effect ( $G \rightarrow X$ )) is reasonable if the pleiotropic path  $G \rightarrow Y$  does not branch off to  $X$ .  
56 However, if there is such a branching off, the variable representing the split is a confounder of  
57 the  $X - Y$  relationship and we fall back on the violation of the second assumption (exchange-  
58 ability), making it the most problematic. Therefore, in this paper, we extend the standard MR  
59 model to incorporate the presence of a latent (i.e. unmeasured) heritable confounder ( $U$ ) and  
60 estimate its contribution to traits  $X$  and  $Y$ , while simultaneously estimating the bi-directional  
61 causal effect between the two traits. Standard MR methods are vulnerable to such heritable  
62 confounders, since any genetic marker directly associated with the confounder may be selected  
63 as an instrument for the exposure. However, such instruments will have a direct effect on the  
64 outcome that is correlated to their instrument strength, violating the InSIDE assumption and  
65 biasing the causal effect estimate.

66 The outline of the paper is as follows: first, the extended MR model is introduced and the  
67 likelihood function for the observed genome-wide summary statistics (for  $X$  and  $Y$ ) is derived.  
68 We then test and compare the method against conventional and more advanced (such as CAUSE

69 [10] and MR RAPS [11]) MR approaches through extensive simulation settings, including several  
70 violations of the model assumptions. Finally, the approach is applied to association summary  
71 statistics (based on the UK Biobank and meta-analysis studies) of 13 complex traits to re-assess  
72 all pairwise bi-directional causal relationships between them.

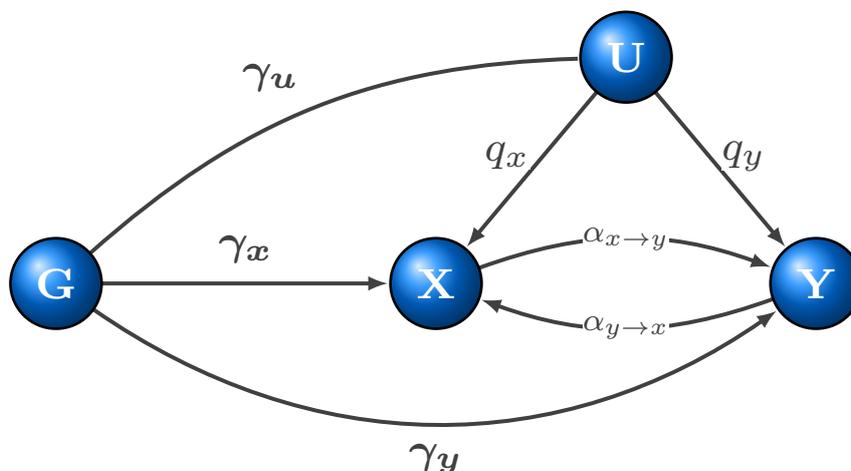
## 73 2 Methods

### 74 2.1 The underlying structural equation model

75 Let  $X$  and  $Y$  denote continuous random variables representing two complex traits. Let us  
76 assume (for simplicity) that there is one heritable confounder  $U$  of these traits. To simplify  
77 notation we assume that  $E(X) = E(Y) = E(U) = 0$  and  $Var(X) = Var(Y) = Var(U) = 1$ .  
78 The genome-wide sequence data for  $M$  sequence variants is denoted by  $G = (G_1, G_2, \dots, G_M)$ .  
79 The aim of our work is to dissect the effects of the heritable confounding factor  $U$  from the  
80 bi-directional causal effects of these two traits ( $X$  and  $Y$ ). For this we consider a model (see  
81 Figure 1) defined by the following equations:

$$\begin{aligned} X &= q_x \cdot U + \alpha_{y \rightarrow x} Y + G \cdot \gamma_x + e_x && \text{with } e_x \sim \mathcal{N}(0, \nu_x^2) \\ Y &= q_y \cdot U + \alpha_{x \rightarrow y} X + G \cdot \gamma_y + e_y && \text{with } e_y \sim \mathcal{N}(0, \nu_y^2) \\ U &= G \cdot \gamma_u + e_u && \text{with } e_u \sim \mathcal{N}(0, \nu_u^2) \end{aligned}$$

82 where  $\gamma_x, \gamma_y, \gamma_u \in \mathcal{R}^M$  denote the (true multivariable) direct effect of all  $M$  genetic variants  
83 on  $X, Y$  and  $U$ , respectively. All error terms ( $e_x, e_y$  and  $e_u$ ) are assumed to be independent of  
84 each other and normally distributed with variances  $\nu_x^2, \nu_y^2$  and  $\nu_u^2$ , respectively.



**Figure 1: Schematic representation of the extended structural equation model (SEM).**  $X$  and  $Y$  are two complex traits under scrutiny with a latent (heritable) confounder  $U$  with causal effects  $q_x$  and  $q_y$  on them.  $G$  represents a genetic instrument, with effects  $\gamma_x, \gamma_y$  and  $\gamma_u$ , respectively. Traits  $X$  and  $Y$  have causal effects on each other, which are denoted by  $\alpha_{x \rightarrow y}$  and  $\alpha_{y \rightarrow x}$ .

85 Note that we do not include in the model reverse causal effects on the confounder ( $X \rightarrow U$  and  
86  $Y \rightarrow U$ ). The reason for this is the following: Let  $s_x$  and  $s_y$  denote those causal effect of  $X$  and  
87  $Y$  on  $U$ . We can see that by reparameterising the original model to  $\alpha'_{x \rightarrow y} := \alpha_{x \rightarrow y} + s_x \cdot q_y$ ,  
88  $\alpha'_{y \rightarrow x} := \alpha_{y \rightarrow x} + s_y \cdot q_x$  and  $q'_x := q_x / (1 - q_x \cdot s_x)$ ,  $q'_y := q_y / (1 - q_y \cdot s_y)$ , the genetic effects produced  
89 by the extended model with reverse causal effects on  $U$  and the simpler model (Figure 1) with  
90 the updated parameters are indistinguishable. Thus these extra parameters are not identifiable

91 and the reparameterisation means that  $\alpha_{x \rightarrow y}$  and  $\alpha_{y \rightarrow x}$  in our model represent the total causal  
92 effects, some of which may be mediated by  $U$ .

93 Note that the model cannot be represented by classical directed acyclic graphs, as the bi-  
94 directional causal effects could form a cycle. However, the equations can be reorganised to  
95 avoid recursive formulation as follows:

$$\begin{aligned} X &= q_x \cdot U + \alpha_{y \rightarrow x} \cdot (q_y \cdot U + \alpha_{x \rightarrow y} X + G \cdot \gamma_y + e_y) + G \cdot \gamma_x + e_x \\ Y &= q_y \cdot U + \alpha_{x \rightarrow y} \cdot (q_x \cdot U + \alpha_{y \rightarrow x} Y + G \cdot \gamma_x + e_x) + G \cdot \gamma_y + e_y \\ U &= G \cdot \gamma_u + e_u \end{aligned}$$

96 Regrouping the terms gives

$$\begin{aligned} (1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}) \cdot X &= (q_x + \alpha_{y \rightarrow x} \cdot q_y) \cdot U + \alpha_{y \rightarrow x} (G \cdot \gamma_y) + G \cdot \gamma_x + (e_x + \alpha_{y \rightarrow x} \cdot e_y) \\ (1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}) \cdot Y &= (q_y + \alpha_{x \rightarrow y} \cdot q_x) \cdot U + \alpha_{x \rightarrow y} (G \cdot \gamma_x) + G \cdot \gamma_y + (e_y + \alpha_{x \rightarrow y} \cdot e_x) \\ U &= G \cdot \gamma_u + e_u \end{aligned}$$

97 Substituting  $U$  into the first two equations yields

$$\begin{aligned} X &= \frac{q_x + \alpha_{y \rightarrow x} \cdot q_y}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot (G \cdot \gamma_u) + \frac{\alpha_{y \rightarrow x}}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} (G \cdot \gamma_y) + \frac{1}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} (G \cdot \gamma_x) + e_x \\ Y &= \frac{q_y + \alpha_{x \rightarrow y} \cdot q_x}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot (G \cdot \gamma_u) + \frac{\alpha_{x \rightarrow y}}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} (G \cdot \gamma_x) + \frac{1}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} (G \cdot \gamma_y) + e_y \end{aligned}$$

98 with

$$\begin{aligned} e_x &:= \frac{e_x + \alpha_{y \rightarrow x} \cdot e_y + (q_x + \alpha_{y \rightarrow x} \cdot q_y) \cdot e_u}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \sim \mathcal{N}(0, i_x) \\ e_y &:= \frac{e_y + \alpha_{x \rightarrow y} \cdot e_x + (q_y + \alpha_{x \rightarrow y} \cdot q_x) \cdot e_u}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \sim \mathcal{N}(0, i_y) \end{aligned}$$

99 where  $i_x := (\nu_x^2 + \alpha_{y \rightarrow x}^2 \nu_y^2 + (q_x + \alpha_{y \rightarrow x} q_y) \nu_u^2) / (1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x})^2$  and  $i_y := (\nu_y^2 + \alpha_{x \rightarrow y}^2 \nu_x^2 +$   
100  $(q_y + \alpha_{x \rightarrow y} q_x) \nu_u^2) / (1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x})^2$ . Note that  $i_x$  is equivalent to the LD score regression  
101 intercept<sup>[12]</sup>.

102 We model the genetic architecture of these direct effects with a spike-and-slab distribution,  
103 assuming that only  $0 \leq \pi_x, \pi_y, \pi_u \leq 1$  proportion of the genome have a direct effect on  $X, Y, U$ ,  
104 respectively and these direct effects come from a Gaussian distribution. Namely,

$$\begin{aligned} \gamma_x &= \zeta_x \odot \kappa_x & \text{with } \kappa_x &\sim \mathcal{N}(0, \sigma_x^2 \cdot I) \text{ and } \zeta_x \sim \mathcal{B}_m(1, \pi_x) \\ \gamma_y &= \zeta_y \odot \kappa_y & \text{with } \kappa_y &\sim \mathcal{N}(0, \sigma_y^2 \cdot I) \text{ and } \zeta_y \sim \mathcal{B}_m(1, \pi_y) \\ \gamma_u &= \zeta_u \odot \kappa_u & \text{with } \kappa_u &\sim \mathcal{N}(0, \sigma_u^2 \cdot I) \text{ and } \zeta_u \sim \mathcal{B}_m(1, \pi_u) \end{aligned}$$

105 Here,  $\odot$  denotes element-wise multiplication and  $\mathcal{B}_m(1, q)$  the  $m$  dimensional independent Bernoulli  
106 distribution. Further, we assume that all  $\kappa_x, \kappa_y, \kappa_u$ s are independent of each other and so are  
107 all  $\zeta_x, \zeta_y, \zeta_u$ s. We can refer to  $h_x^2 := M \cdot \pi_x \cdot \sigma_x^2$  as the *direct heritability* of  $X$ , i.e. independent  
108 of the genetic basis of  $U$  and  $Y$ . Similar notation is adapted for  $U$  ( $h_u^2 := M \cdot \pi_u \cdot \sigma_u^2$ ) and  $Y$   
109 ( $h_y^2 := M \cdot \pi_y \cdot \sigma_y^2$ ). Note that when  $q_x = 0$  and  $q_y \neq 0$  (or vice versa), this means that there is  
110 no confounder  $U$  present, but the genetic architecture of  $Y$  (or  $X$ ) can be better described by a  
111 three component Gaussian mixture distribution.

112 We assume that the correlation (across markers) between the direct effects of a genetic variant  
113 on  $X, Y$  and  $U$  is zero, i.e.  $cov(\gamma_x, \gamma_y) = cov(\gamma_x, \gamma_u) = cov(\gamma_u, \gamma_y) = 0$ . Note that this

114 assumption still allows for a potential correlation between the total effect of  $G$  on  $X$  and its  
115 horizontal pleiotropic effect on  $Y$ , but only due to the confounder  $U$  and through the reverse  
116 causal effect  $Y \rightarrow X$ . As we argued above, this is a reasonable assumption, since the most plau-  
117 sible reason (apart from outcome-dependent sampling, which is out of the scope of this paper)  
118 for the violation of the InSIDE assumption may be one or more heritable confounder(s).

119 For simplicity, we also assume that the set of genetic variants with direct effects on each trait  
120 overlap only randomly, i.e. the fraction of the genome directly associated with both  $X$  and  $Y$   
121 is  $\pi_x \cdot \pi_y$ , etc. This assumption is in line with recent observation that the bulk of observed  
122 pleiotropy can be explained by extreme polygenicity with random overlap between trait loci<sup>[13]</sup>.  
123 Note that uncorrelated effects (e.g.  $\text{cov}(\gamma_x, \gamma_y) = 0$ ) do not ensure that the active variant sets  
124 overlap randomly, this is a slightly stronger assumption.

## 125 2.2 The observed association summary statistics

126 Let us now assume that we observe univariable association summary statistics for these two  
127 traits from two (potentially overlapping) finite samples  $N_x$  and  $N_y$  of size  $n_x, n_y$ , respectively.  
128 In the following, we will derive observed summary statistics in sample  $N_x$  and then we will  
129 repeat the analogous exercise for sample  $N_y$ . Let the realisations of  $X, Y$  and  $U$  be denoted by  
130  $\mathbf{x}, \mathbf{y}$  and  $\mathbf{u} \in \mathcal{R}^{n_x}$ . The genome-wide genetic data is represented by  $\mathbf{G}_x \in \mathcal{R}^{n_x \times M}$  and the genetic  
131 data for a single nucleotide polymorphism (SNP)  $k$  tested for association is  $\mathbf{g}_k \in \mathcal{R}^{n_x}$ . Note the  
132 distinction between the  $k$ -th column of  $\mathbf{G}_x$ , which is the  $k$ -th sequence variant, in contrast to  $\mathbf{g}_k$ ,  
133 which is the  $k$ -th SNP tested for association in the GWAS. We assume that all SNP genotypes  
134 have been standardised to have zero mean and unit variance. The marginal effect size estimate  
135 for SNP  $k$  of trait  $X$  can then be written as  $\hat{\beta}_k^x = \mathbf{g}_k' \cdot \mathbf{x} / n_x$ , which is a special case of univariable  
136 standard normal linear regression when both the outcome and the predictor is standardised to  
137 have zero mean and unit variance<sup>[12]</sup>. Note that  $\mathbf{x}'$  denotes the transposed of the column vector  
138  $\mathbf{x}$ . This can be further transformed as

$$\begin{aligned} \hat{\beta}_k^x &= \mathbf{g}_k' \cdot \mathbf{x} / n_x \\ &= \frac{q_x + \alpha_{y \rightarrow x} \cdot q_y}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \mathbf{g}_k' \cdot \mathbf{G}_x \cdot \boldsymbol{\gamma}_u / n_x + \frac{\alpha_{y \rightarrow x}}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \mathbf{g}_k' \cdot \mathbf{G}_x \cdot \boldsymbol{\gamma}_y / n_x \\ &\quad + \frac{1}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \mathbf{g}_k' \cdot \mathbf{G}_x \cdot \boldsymbol{\gamma}_x / n_x + \mathbf{g}_k' \cdot \boldsymbol{\epsilon}_x / n_x \end{aligned}$$

139 By denoting the linkage disequilibrium (LD) between variant  $k$  and all markers in the genome  
140 with  $\boldsymbol{\rho}_k = \mathbf{G}_x' \cdot \mathbf{g}_k / n_x$  we get

$$\hat{\beta}_k^x = \frac{q_x + \alpha_{y \rightarrow x} \cdot q_y}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \boldsymbol{\rho}_k' \cdot \boldsymbol{\gamma}_u + \frac{\alpha_{y \rightarrow x}}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \boldsymbol{\rho}_k' \cdot \boldsymbol{\gamma}_y + \frac{1}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \boldsymbol{\rho}_k' \cdot \boldsymbol{\gamma}_x + \eta_k^x$$

141 with  $\eta_k^x := \mathbf{g}_k' \cdot \boldsymbol{\epsilon}_x / n_x \sim \mathcal{N}(0, i_x / n_x)$ . Given the above-defined genetic effect size distribution the  
142 equation becomes

$$\begin{aligned} \hat{\beta}_k^x &= \frac{q_x + \alpha_{y \rightarrow x} \cdot q_y}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \underbrace{\boldsymbol{\rho}_k' \cdot (\boldsymbol{\zeta}_u \odot \boldsymbol{\kappa}_u)}_{z_k^{(u)}} + \frac{\alpha_{y \rightarrow x}}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \underbrace{\boldsymbol{\rho}_k' \cdot (\boldsymbol{\zeta}_y \odot \boldsymbol{\kappa}_y)}_{z_k^{(y)}} \\ &\quad + \frac{1}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot \underbrace{\boldsymbol{\rho}_k' \cdot (\boldsymbol{\zeta}_x \odot \boldsymbol{\kappa}_x)}_{z_k^{(x)}} + \eta_k^x \\ &= \frac{q_x + \alpha_{y \rightarrow x} \cdot q_y}{1 - \alpha_{y \rightarrow x} \alpha_{y \rightarrow x}} \cdot z_k^{(u)} + \frac{\alpha_{y \rightarrow x}}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \cdot z_k^{(y)} + \frac{1}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} z_k^{(x)} + \eta_k^x \end{aligned}$$

143 Similarly, assuming that the LD structures ( $\rho_k$ ) in the two samples are comparable, for  $\widehat{\beta}_k^y$   
144 estimated in the other sample ( $N_y$ ) we obtain

$$\widehat{\beta}_k^y = \frac{\alpha_{x \rightarrow y} \cdot q_x + q_y}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \cdot z_k^{(u)} + \frac{\alpha_{x \rightarrow y}}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \cdot z_k^{(x)} + \frac{1}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} z_k^{(y)} + \eta_k^y$$

145 with  $\eta_k^y \sim \mathcal{N}(0, i_y/n_y)$ .

146 Therefore, the joint effect size estimates can be written as

$$\begin{pmatrix} \widehat{\beta}_k^x \\ \widehat{\beta}_k^y \end{pmatrix} = \frac{1}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \left( \begin{pmatrix} \alpha_{y \rightarrow x} \cdot q_y + q_x \\ \alpha_{x \rightarrow y} \cdot q_x + q_y \end{pmatrix} z_k^{(u)} + \begin{pmatrix} 1 \\ \alpha_{x \rightarrow y} \end{pmatrix} z_k^{(x)} + \begin{pmatrix} \alpha_{y \rightarrow x} \\ 1 \end{pmatrix} z_k^{(y)} \right) + \begin{pmatrix} \eta_k^x \\ \eta_k^y \end{pmatrix}$$

147 Following the same rationale as the cross-trait LD score regression<sup>[14]</sup>, the noise term distribution  
148 is readily obtained

$$\begin{pmatrix} \eta_k^x \\ \eta_k^y \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} i_x/n_x & \frac{n_{x \cap y}}{n_x \cdot n_y} \cdot r_{x,y} \\ \frac{n_{x \cap y}}{n_x \cdot n_y} \cdot r_{x,y} & i_y/n_y \end{pmatrix} \right)$$

149 where  $r_{x,y}$  is the observational correlation between variables  $X$  and  $Y$  and  $n_{x \cap y}$  is the size of  
150 the overlapping samples for  $X$  and  $Y$ . Since both  $n_{x \cap y}$  and  $r_{x,y}$  cannot be estimated, we simply  
151 denote  $i_{x,y} := r_{x,y} \cdot \frac{n_{x \cap y}}{\sqrt{n_x \cdot n_y}}$  as the only estimated parameter and parameterise the covariance  
152 term as  $\frac{i_{x,y}}{\sqrt{n_x \cdot n_y}}$ . Note that  $i_{x,y}$  is the cross-trait LD score regression intercept.

153 The bivariate probability density function (PDF) of these summary statistics cannot be obtained  
154 analytically, but in the following we demonstrate that the characteristic function can be derived.  
155 Let us first compute the characteristic function of this two-dimensional random variable, know-  
156 ing that  $z_k^{(x)}, z_k^{(u)}, z_k^{(y)}$  and  $(\eta_k^x, \eta_k^y)$  are independent, hence the characteristic function can be  
157 factorised:

$$\begin{aligned} \varphi_{(\widehat{\beta}_k^x, \widehat{\beta}_k^y)}(v, w) &= E \left[ \exp \left( i \cdot (v \cdot \widehat{\beta}_k^x + w \cdot \widehat{\beta}_k^y) \right) \right] \\ &= E \left[ \exp \left( i \cdot \left( v \cdot \left( \frac{z_k^{(x)} + (\alpha_{y \rightarrow x} \cdot q_y + q_x) \cdot z_k^{(u)} + \alpha_{y \rightarrow x} \cdot z_k^{(y)}}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} + \eta_k^x \right) + \right. \right. \right. \\ &\quad \left. \left. \left. + w \cdot \left( \frac{z_k^{(y)} + (\alpha_{x \rightarrow y} \cdot q_x + q_y) \cdot z_k^{(u)} + \alpha_{x \rightarrow y} \cdot z_k^{(x)}}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} + \eta_k^y \right) \right) \right) \right] \\ &= E \left[ \exp \left( i \cdot z_k^{(u)} \cdot \frac{v \cdot (\alpha_{y \rightarrow x} \cdot q_y + q_x) + w \cdot (\alpha_{x \rightarrow y} \cdot q_x + q_y)}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \right] \\ &\quad \times E \left[ \exp \left( i \cdot z_k^{(x)} \cdot \frac{v + \alpha_{x \rightarrow y} \cdot w}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \right] \cdot E \left[ \exp \left( i \cdot z_k^{(y)} \cdot \frac{w + \alpha_{y \rightarrow x} \cdot v}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \right] \\ &\quad \times E \left[ \exp \left( i \cdot (v \cdot \eta_k^x + w \cdot \eta_k^y) \right) \right] \\ &= \varphi_{z_k^{(u)}} \left( \frac{v \cdot (\alpha_{y \rightarrow x} \cdot q_y + q_x) + w \cdot (\alpha_{x \rightarrow y} \cdot q_x + q_y)}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \\ &\quad \times \varphi_{z_k^{(x)}} \left( \frac{v + \alpha_{x \rightarrow y} \cdot w}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \cdot \varphi_{z_k^{(y)}} \left( \frac{w + \alpha_{y \rightarrow x} \cdot v}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \cdot \varphi_{(\eta_k^x, \eta_k^y)}(v, w) \end{aligned}$$

158 In the following we will work out each of the characteristic functions on the right hand side.

159 **2.3 Characteristic functions of  $z_k^{(u)}$ ,  $z_k^{(x)}$ ,  $z_k^{(y)}$  and  $(\eta_k^x, \eta_k^y)$**

160 It is reasonable to assume that linkage disequilibrium (LD) fades off beyond 1Mb distance. Thus,  
161 without loss of generality we can assume that non-zero LD does not extend beyond  $m_0$  markers  
162 around the focal variant. Hence we can assume that the length of  $\rho_k$  is  $m_0$  and only consider  
163  $\gamma_x, \gamma_y$  and  $\gamma_u$  to be of length  $m_0$  instead of  $m$ . Let us first approximate the distribution of  $\rho_k$   
164 values following a spike and slab Gaussian mixture, i.e. proportion  $\pi_k$  of the  $m_0$  SNPs have  
165 non-zero LD, coming from a Gaussian distribution  $\mathcal{N}(0, \sigma_k^2)$  and the remaining  $(1 - \pi_k)$  fraction  
166 of the LD values is zero. In mathematical notation

$$\rho_k = \mathbf{r}_k \odot \boldsymbol{\kappa}_k \quad \text{with} \quad \mathbf{r}_k \sim \mathcal{N}(0, \sigma_k^2 \cdot I) \quad \text{and} \quad \boldsymbol{\kappa}_k \sim \mathcal{B}_{m_0}(1, \pi_k)$$

167 Therefore  $z_k^{(u)}$  can be written of the form

$$\begin{aligned} z_k^{(u)} &= \boldsymbol{\rho}_k' \cdot (\boldsymbol{\zeta}_u \odot \boldsymbol{\kappa}_u) = (\mathbf{r}_k \odot \boldsymbol{\kappa}_k)' \cdot (\boldsymbol{\zeta}_u \odot \boldsymbol{\kappa}_u) = (\mathbf{r}_k \odot \boldsymbol{\zeta}_u)' \cdot \underbrace{(\boldsymbol{\kappa}_k \odot \boldsymbol{\kappa}_u)}_{\boldsymbol{\kappa}_{k,u} \sim \mathcal{B}(1, \pi_k \cdot \pi_u)} \\ &= \sum_{j=1}^{m_0} (\mathbf{r}_k \odot \boldsymbol{\zeta}_u)_j \cdot (\boldsymbol{\kappa}_{k,u})_j \end{aligned}$$

168 The PDF of the product of two zero-mean Gaussians ( $\mathbf{r}_k$  and  $\boldsymbol{\zeta}_u$ ) is a modified Bessel function  
169 of the second kind of order zero ( $K_0(\omega)$ ) [15], more precisely

$$f_{(\mathbf{r}_k \odot \boldsymbol{\zeta}_u)_j}(t) = \frac{1/\pi}{\sigma_u \cdot \sigma_k} \cdot K_0\left(\frac{|t|}{\sigma_u \cdot \sigma_k}\right)$$

170 and its characteristic function [16,17] is

$$\varphi_0(t) = E(\exp(i \cdot t \cdot (\mathbf{r}_k \odot \boldsymbol{\zeta}_u)_j)) = \frac{1}{\sqrt{\sigma_u^2 \cdot \sigma_k^2 \cdot t^2 + 1}}$$

171 Next, the characteristic function of the product of  $(\mathbf{r}_k \odot \boldsymbol{\zeta}_u)_j$  and a Bernoulli distributed  $(\boldsymbol{\kappa}_{k,u})_j$   
172 is

$$\begin{aligned} \varphi_1(t) &= E(\exp(i \cdot t \cdot (\mathbf{r}_k \odot \boldsymbol{\zeta}_u)_j) \cdot (\boldsymbol{\kappa}_{k,u})_j) \\ &= \pi_k \cdot \pi_u \cdot E(\exp(i \cdot t \cdot (\mathbf{r}_k \odot \boldsymbol{\zeta}_u)_j)) + (1 - \pi_k \cdot \pi_u) \cdot E(\exp(i \cdot t \cdot 0)) \\ &= \pi_k \cdot \pi_u \cdot \varphi_0(t) + (1 - \pi_k \cdot \pi_u) \\ &= \frac{\pi_k \cdot \pi_u}{\sqrt{\sigma_u^2 \cdot \sigma_k^2 \cdot t^2 + 1}} + (1 - \pi_k \cdot \pi_u) \end{aligned}$$

173 Hence the characteristic function of the sum of  $m_0$  independent random variables is the product  
174 of them, we have

$$\varphi_{z_k^{(u)}}(t) = \left( \frac{\pi_k \cdot \pi_u}{\sqrt{\sigma_u^2 \cdot \sigma_k^2 \cdot t^2 + 1}} + (1 - \pi_k \cdot \pi_u) \right)^{m_0}$$

175 Finally, we apply a first order Taylor series approximation (around 1) of the log of the charac-  
176 teristic function in order to speed up computation and improve numerical accuracy

$$\begin{aligned} \log(\varphi_{z_k^{(u)}}(t)) &= m_0 \cdot \log \left( \frac{\pi_k \cdot \pi_u}{\sqrt{\sigma_u^2 \cdot \sigma_k^2 \cdot t^2 + 1}} + (1 - \pi_k \cdot \pi_u) \right) \\ &= m_0 \cdot \log \left( 1 - \pi_k \cdot \pi_u \cdot \left( 1 - \frac{1}{\sqrt{\sigma_u^2 \cdot \sigma_k^2 \cdot t^2 + 1}} \right) \right) \\ &\approx -m_0 \cdot \pi_k \cdot \pi_u \cdot \left( 1 - \frac{1}{\sqrt{\sigma_u^2 \cdot \sigma_k^2 \cdot t^2 + 1}} \right) \end{aligned}$$

177 Analogously, the approximation of the logarithm of the characteristic functions of  $z_k^{(x)}$  and  $z_k^{(y)}$   
178 is

$$\begin{aligned} \log(\varphi_{z_k^{(x)}}(t)) &\approx -m_0 \cdot \pi_k \cdot \pi_x \cdot \left( 1 - \frac{1}{\sqrt{\sigma_x^2 \cdot \sigma_k^2 \cdot t^2 + 1}} \right) \\ \log(\varphi_{z_k^{(y)}}(t)) &\approx -m_0 \cdot \pi_k \cdot \pi_y \cdot \left( 1 - \frac{1}{\sqrt{\sigma_y^2 \cdot \sigma_k^2 \cdot t^2 + 1}} \right) \end{aligned}$$

179 Since the characteristic function of a centred multivariate Gaussian with variance-covariance  
180 matrix  $\Sigma$  is  $\exp(-(1/2) \cdot \mathbf{t}' \cdot \Sigma \cdot \mathbf{t})$  we have

$$\log \left( \varphi_{(\eta_k^x, \eta_k^y)}(v, w) \right) = -\frac{1}{2} \cdot \left( \frac{i_x}{n_x} \cdot v^2 + 2 \cdot \frac{i_{x,y}}{\sqrt{n_x \cdot n_y}} \cdot v \cdot w + \frac{i_y}{n_y} \cdot w^2 \right)$$

## 181 2.4 From characteristic function to probability density function

182 The final form of the logarithm of the joint characteristic function of the transformed summary  
183 statistics is

$$\begin{aligned}
 \log \left( \varphi_{(\hat{\beta}_k^x, \hat{\beta}_k^y)}(v, w) \right) &= \log \left( \varphi_{z_k^{(x)}} \left( \frac{v + \alpha_{x \rightarrow y} w}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \right) + \log \left( \varphi_{z_k^{(y)}} \left( \frac{w + \alpha_{y \rightarrow x} v}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \right) \\
 &+ \log \left( \varphi_{z_k^{(u)}} \left( \frac{v \cdot (\alpha_{y \rightarrow x} \cdot q_y + q_x) + w \cdot (\alpha_{x \rightarrow y} \cdot q_x + q_y)}{1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x}} \right) \right) \\
 &+ \log \left( \varphi_{(\eta_k^x, \eta_k^y)}(v, w) \right) \\
 &\approx -m_0 \cdot \pi_k \cdot \pi_x \cdot \left( 1 - \frac{1}{\sqrt{\frac{\sigma_x^2 \cdot \sigma_k^2 \cdot (v + \alpha_{x \rightarrow y} w)^2}{(1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x})^2} + 1}} \right) \\
 &- m_0 \cdot \pi_k \cdot \pi_y \cdot \left( 1 - \frac{1}{\sqrt{\frac{\sigma_y^2 \cdot \sigma_k^2 \cdot (w + \alpha_{y \rightarrow x} v)^2}{(1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x})^2} + 1}} \right) \\
 &- m_0 \cdot \pi_k \cdot \pi_u \cdot \left( 1 - \frac{1}{\sqrt{\frac{\sigma_u^2 \cdot \sigma_k^2 \cdot (v \cdot (\alpha_{y \rightarrow x} \cdot q_y + q_x) + w \cdot (\alpha_{x \rightarrow y} \cdot q_x + q_y))^2}{(1 - \alpha_{x \rightarrow y} \alpha_{y \rightarrow x})^2} + 1}} \right) \\
 &- \frac{1}{2} \cdot \left( \frac{i_x}{n_x} \cdot v^2 + 2 \cdot \frac{i_{x,y}}{\sqrt{n_x \cdot n_y}} \cdot v \cdot w + \frac{i_y}{n_y} \cdot w^2 \right)
 \end{aligned} \tag{1}$$

184 Using the inversion theorem for characteristic functions we can express the joint distribution of  
185  $(\hat{\beta}_k^x, \hat{\beta}_k^y)$  as

$$f_{(\hat{\beta}_k^x, \hat{\beta}_k^y)}(x, y) = \left( \frac{1}{2\pi} \right)^2 \cdot \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-i \cdot (x \cdot v + y \cdot w)) \cdot \varphi_{(\hat{\beta}_k^x, \hat{\beta}_k^y)}(v, w) \, dv \, dw$$

186 This integral can be efficiently computed by Fast Fourier Transformation (FFT, see [18](#) and  
187 references within). To speed up computation, we bin SNPs according to their  $\pi_k$  and  $\sigma_k$  values  
188 ( $10 \times 10$  bins with equidistant centres) and for SNPs in the same bin the PDF function is  
189 evaluated over a fine grid ( $2^7 \times 2^7$  combinations) using the FFT.

190 To reduce the number of parameters we define  $t_x := \sigma_u \cdot q_x$  and  $t_y := \sigma_u \cdot q_y$  since  $\sigma_u$  and  $q_x$  are  
191 separately not identifiable, but only their product is. Similarly  $\pi_u$  is unidentifiable, and is set to  
192 an arbitrary value of 0.1. For improved interpretability, we slightly reparameterise the likelihood  
193 function by using  $h_x^2 := \pi_x \cdot M \cdot \sigma_x^2$ ,  $h_y^2 := \pi_y \cdot M \cdot \sigma_y^2$ . Since different SNPs are correlated we have to  
194 estimate the over-counting of each SNP. We choose the same strategy as LD score regression [12](#)  
195 and weigh each SNP by the inverse of its restricted LD score, i.e.  $w_k = 1 / \sum_{j=1}^{m_0} r_{jk}^2$ , where  $r_{jk}$   
196 is the correlation between GWAS SNPs  $k$  and  $j$ . The log-likelihood function is, thus, of the  
197 form

$$\log \left( \mathcal{L} \left( \boldsymbol{\theta} \left| \begin{pmatrix} \hat{\beta}^x \\ \hat{\beta}^y \end{pmatrix} \right. \right) \right) \propto \sum_{k=1}^K w_k \cdot f_k \left( \hat{\beta}_k^x, \hat{\beta}_k^y \right) \tag{2}$$

198 where  $f_k \left( \hat{\beta}_k^x, \hat{\beta}_k^y \right)$  is the log-likelihood function value for SNP  $k$ . Parameters  $\{n_x, n_y, m, \sigma_{k=1, \dots, K}, \pi_{k=1, \dots, K}\}$   
199 are known and the other 11 parameters

$$\boldsymbol{\theta} = \{ \pi_x, \pi_y, h_x^2, h_y^2, t_x, t_y, \alpha_{x \rightarrow y}, \alpha_{y \rightarrow x}, i_x, i_y, i_{x,y} \}$$

200 are to be estimated from the observed association summary statistics. In order to speed up  
 201 computation, we can estimate the 11 parameters in two separate steps: the first estimates for  
 202 each trait the parameters  $\pi_x, i_x$  and  $\pi_y, i_y$  (SNP polygenicity and LD-score intercept) and the  
 203 total heritability (unlike the direct heritability obtained by the full-model of LHC-MR) by using  
 204 a simplified model with only the trait of interest, without a second trait or confounder, e.g.  
 205 we fit only  $\pi_x, h_x^2$  and  $i_x$  using  $\hat{\beta}^x$  and assume that  $\pi_x$  and  $i_x$  do not change when two traits  
 206 are taken into account. Note that  $\pi_x$  may change slightly (decreasing from the total- to direct  
 207 polygenicity), but its value has little impact on the likelihood function. The estimates from the  
 208 first step can then be fixed for the parameter estimation of trait pairs. Since only  $\pi_x, i_x$  and  
 209  $\pi_y, i_y$  are fixed, the remaining parameters to estimate are now:

$$\theta = \{h_x^2, h_y^2, t_x, t_y, \alpha_{x \rightarrow y}, \alpha_{y \rightarrow x}, i_{x,y}\}$$

210 It is key to note that our approach does not aim to estimate individual (direct or indirect) SNP  
 211 effects, as these are handled as random effects. By replacing  $U$  with  $-U$  we swap the signs of  
 212 both  $t_x$  and  $t_y$ , therefore these parameters are unique only if the sign of one of them is fixed.  
 213 Thus, we will have the following restrictions on the parameter ranges:  $h_x^2, h_y^2, t_x$  are in  $[0, 1]$ ,  
 214  $t_y, \alpha_{x \rightarrow y}, \alpha_{y \rightarrow x}, i_{x,y}$  are in  $[-1, 1]$ .

## 215 2.5 Likelihood maximisation and standard error calculation

216 Our method, termed *Latent Heritable Confounder Mendelian Randomisation (LHC-MR)*, max-  
 217 imises this likelihood function to obtain the maximum likelihood estimate (MLE). Due to the  
 218 complexity of the likelihood surface, we initialise the maximisation using 50 different starting  
 219 points, where they come from a uniform distribution within the parameter-specific ranges men-  
 220 tioned above. We then choose parameter estimates corresponding to the highest likelihood of the  
 221 50 runs. Run time depends on the number of iterations during the maximisation procedure, and  
 222 is linear with respect to the number of SNPs. It takes  $\sim 0.25$  CPU-minute to fit the complete  
 223 model to 50,000 SNPs with a single starting point.

224 Given the particular nature of the underlying directed graph, two different sets of parameters  
 225 lead to an identical fit of the data, resulting in two global optima. The reason for this is  
 226 the difficulty in distinguishing the ratio of the confounder effects ( $t_y/t_x$ ) from the causal effect  
 227 ( $\alpha_{x \rightarrow y}$ ), as illustrated in Figure [S2](#) by the slopes belonging to different SNP-clusters. More  
 228 rigorously, it can be show that if  $\{h_x, h_y, \alpha_{x \rightarrow y}, \alpha_{y \rightarrow x}, t_x, t_y\}$  is an optimum, then so will be  
 229  $\{h'_x, h'_y, \alpha'_{x \rightarrow y}, \alpha'_{y \rightarrow x}, t'_x, t'_y\}$ , where

$$\begin{aligned} h'_x &= t_x + t_y \cdot \alpha_{y \rightarrow x} \\ h'_y &= h_y \\ \alpha'_{x \rightarrow y} &= \frac{\alpha_{x \rightarrow y} + w}{1 + \alpha_{y \rightarrow x} \cdot w} \\ \alpha'_{y \rightarrow x} &= \alpha_{y \rightarrow x} \\ t'_x &= h_x \cdot (1 + \alpha_{y \rightarrow x} \cdot w) \\ t'_y &= -h_x \cdot w \end{aligned}$$

230 with  $w = t_y/t_x$  (for further derivations, see Supplementary Section [1.1](#)). This allows us to  
 231 directly obtain both optima, even if the optimisation only revealed one of them. It happens  
 232 very often that one of these parameter sets are outside of the allowed ranges and hence can  
 233 be automatically excluded. If not, we keep track of both parameter estimates maximising the  
 234 likelihood function. Note that, we call the one for which the direct heritability is larger than  
 235 the indirect one, i.e.  $h_x^2 > t_x^2$ , the primary solution. We show that for real data application this

236 solution is far more plausible than the alternative optimum. Finally, note that such bimodality  
237 can be observed at different levels: (i) For one given data generation, using multiple starting  
238 points leads to different optima; (ii) LHC-MR applied to multiple different data generations for  
239 a fixed parameter setting can yield different optima. Both of these situations are signs of the  
240 same underlying phenomenon and most often co-occur.

241 We implemented the block jackknife procedure that is also used by LD score regression to  
242 calculate the standard errors. For this we split the genome into 200 jackknife blocks and compute  
243 MLE in a leave-one-block-out fashion yielding  $\hat{\theta}^{(-i)}, i = 1, \dots, 200$  estimates. The variance  
244 of the full SNP MLE is then defined as  $Var(\hat{\theta}) := \frac{m-m \cdot (1/200)}{m \cdot (1/200)} \cdot \frac{1}{200-1} \sum_{i=1}^{200} (\hat{\theta}^{(-i)} - \hat{\theta})^2 =$   
245  $\sum_{i=1}^{200} (\hat{\theta}^{(-i)} - \hat{\theta})^2$ .

## 246 2.6 Decomposition of genetic correlation

247 Given the starting equations for  $X$  and  $Y$  we can calculate their genetic correlation. Denoting  
248 the total (multivariate) genetic effect for  $X$  and  $Y$  as  $\delta_x$  and  $\delta_y$ , we can express them as  
249 follows

$$\begin{aligned}\delta_x &= q_x \cdot \gamma_u + \alpha_{y \rightarrow x} \delta_y + \gamma_x \\ \delta_y &= q_y \cdot \gamma_u + \alpha_{x \rightarrow y} \delta_x + \gamma_y\end{aligned}$$

250 Substituting the second equation to the first yields

$$\begin{aligned}\delta_x &= q_x \cdot \gamma_u + \alpha_{y \rightarrow x} (q_y \cdot \gamma_u + \alpha_{x \rightarrow y} \delta_x + \gamma_y) + \gamma_x \\ &= (q_x + \alpha_{y \rightarrow x} q_y) \cdot \gamma_u + (\alpha_{y \rightarrow x} \alpha_{x \rightarrow y}) \delta_x + \alpha_{y \rightarrow x} \gamma_y + \gamma_x \\ &= ((q_x + \alpha_{y \rightarrow x} q_y) \cdot \gamma_u + \alpha_{y \rightarrow x} \gamma_y + \gamma_x) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})\end{aligned}$$

251 Similarly,

$$\delta_y = ((q_y + \alpha_{x \rightarrow y} q_x) \cdot \gamma_u + \alpha_{x \rightarrow y} \gamma_x + \gamma_y) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})$$

252 Thus the genetic covariance is

$$\begin{aligned}E[\delta_x \cdot \delta_y] &= ((q_x + \alpha_{y \rightarrow x} q_y) \cdot \gamma_u + \alpha_{y \rightarrow x} \gamma_y + \gamma_x) ((q_y + \alpha_{x \rightarrow y} q_x) \cdot \gamma_u + \alpha_{x \rightarrow y} \gamma_x + \gamma_y) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})^2 \\ &= ((q_x + \alpha_{y \rightarrow x} q_y)(q_y + \alpha_{x \rightarrow y} q_x) h_u^2 + \alpha_{y \rightarrow x} h_y^2 + \alpha_{x \rightarrow y} h_x^2) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})^2 \\ &= ((t_x + \alpha_{y \rightarrow x} t_y)(t_y + \alpha_{x \rightarrow y} t_x) + \alpha_{y \rightarrow x} h_y^2 + \alpha_{x \rightarrow y} h_x^2) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})^2\end{aligned}$$

253 and the heritabilities are

$$\begin{aligned}E[\delta_x^2] &= ((t_x + \alpha_{y \rightarrow x} t_y)^2 + \alpha_{y \rightarrow x}^2 h_y^2 + h_x^2) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})^2 \\ E[\delta_y^2] &= ((t_y + \alpha_{x \rightarrow y} t_x)^2 + \alpha_{x \rightarrow y}^2 h_x^2 + h_y^2) / (1 - \alpha_{y \rightarrow x} \alpha_{x \rightarrow y})^2\end{aligned}$$

254 Therefore the genetic correlation takes the form

$$corr(\delta_x, \delta_y) = \frac{(t_x + \alpha_{y \rightarrow x} t_y)(t_y + \alpha_{x \rightarrow y} t_x) + \alpha_{y \rightarrow x} h_y^2 + \alpha_{x \rightarrow y} h_x^2}{\sqrt{((t_x + \alpha_{y \rightarrow x} t_y)^2 + \alpha_{y \rightarrow x}^2 h_y^2 + h_x^2) ((t_y + \alpha_{x \rightarrow y} t_x)^2 + \alpha_{x \rightarrow y}^2 h_x^2 + h_y^2)}} \quad (3)$$

255 These values can be compared to those obtained by LD score regression.

## 256 2.7 Computation of the LD scores

257 We first took 4,773,627 SNPs with info (imputation certainty measure)  $\geq 0.99$  present in the  
258 association summary files from the second round of GWAS by the Neale lab<sup>[19]</sup>. This set was  
259 restricted to 4,650,107 common, high-quality SNPs, defined as being present in both UK10K  
260 and UK Biobank, having MAF  $> 1\%$  in both data sets, non-significant ( $P_{diff} > 0.05$ ) allele  
261 frequency difference between UK Biobank and UK10K and residing outside the HLA region  
262 (chr6:28.5-33.5Mb). For these SNPs, LD scores and regression weights were computed based on  
263 3,781 individuals from the UK10K study<sup>[20]</sup>. To estimate the local LD distribution for each SNP  
264 ( $k$ ), characterised by  $\pi_k, \sigma_k^2$ , we fitted a two-component Gaussian mixture distribution to the  
265 observed local correlations (focal SNP  $\pm 2'500$  markers with MAF  $\geq 0.5\%$  in the UK10K):  
266 (1) one Gaussian component corresponding to zero correlations, reflecting only measurement  
267 noise (whose variance is proportional to the inverse of the reference panel size) and (2) a sec-  
268 ond component with zero mean and a larger variance than the first component (encompassing  
269 measurement noise plus non-zero LD).

## 270 2.8 Simulation settings

271 First, we tested LHC-MR using realistic parameter settings with a mild violation of the classical  
272 MR assumptions. These standard parameter settings consisted of simulating  $m = 234,000$  SNPs  
273 for two non-overlapping cohorts of equal size (for simplicity) of  $n_x = n_y = 50,000$  for each trait.  
274  $X, Y$  and  $U$  were simulated with moderate polygenicity ( $\pi_x = 5 \times 10^{-3}, \pi_y = 1 \times 10^{-2}, \pi_u = 5 \times$   
275  $10^{-2}$ ), and considerable direct heritability ( $h_x^2 = 0.25, h_y^2 = 0.2, h_u^2 = 0.3$ ).  $U$  had a confounding  
276 effect on the two traits as such,  $q_x = 0.3, q_y = 0.2$  (resulting in  $t_x = 0.16, t_y = 0.11$ ), and  $X$  had  
277 a direct causal effect on  $Y$  ( $\alpha_{x \rightarrow y} = 0.3$ ), while the reverse causal effect from  $Y$  to  $X$  was set to  
278 null. Note that in this setting the total heritability of each of these traits is principally driven  
279 by direct effects and less than 10% of the total heritability is through a confounder and in case  
280 of  $Y$  less than an additional 8% of its total heritability is through  $X$ .  
281 It is important to note that for each tested parameter setting, we generated 50 different data  
282 sets, and each data generation underwent a likelihood maximisation of Eq. 2 using 50 starting  
283 points, and produced estimated parameters corresponding to the highest likelihood (simplified  
284 schema in Figure S3).

285 In the following simulations, we changed various parameters of these standard settings to test  
286 the robustness of the method. We explored how increased sample size ( $n_x = n_y = 500,000$ )  
287 or differences in sample sizes ( $(n_x, n_y) = (50,000, 500,000)$  and  $(n_x, n_y) = (500,000, 50,000)$ )  
288 influence causal effect estimates of LHC-MR and other MR methods. We also simulated data  
289 with no causal effect (or with no confounder) and then examined how LHC-MR estimates those  
290 parameters. Next, we varied our causal effects between the two traits by lowering  $\alpha_{x \rightarrow y}$  to 0.1,  
291 and in another setting by introducing a reverse causal effect ( $\alpha_{y \rightarrow x} = -0.1$ ). In addition, we  
292 tried to create extremely unfavourable conditions for all MR analyses by varying the confounding  
293 effects. We did this in several ways: (i) increasing  $q_x$  and  $q_y$  ( $q_x = 0.75, q_y = 0.50$ ), (ii) having  
294 a confounder with causal effects of opposite signs on  $X$  and  $Y$  ( $q_x = 0.3, q_y = -0.2$ ). We also  
295 drastically increased the proportion of SNPs with non-zero effect on traits  $X, Y$  and  $U$  ( $\pi_x, \pi_y$   
296 and  $\pi_u = 0.1, 0.15, 0.2$  respectively). We also simulated data whereby the confounder has lower  
297 ( $\pi_u = 0.01$ ) polygenicity than the two focal traits.

298 Finally, we explored various violations of the assumptions of our model (see Section 2). First,  
299 we introduced two confounders in the simulated data, once with causal effects on  $X$  and  $Y$   
300 that were concordant ( $t_x^{(1)} = 0.16, t_y^{(1)} = 0.11, t_x^{(2)} = 0.22, t_y^{(2)} = 0.16$ ) in sign, and another with  
301 discordant effects ( $t_x^{(1)} = 0.16, t_y^{(1)} = 0.11, t_x^{(2)} = 0.22, t_y^{(2)} = -0.16$ ), while still fitting the model  
302 with only one  $U$ . Second, we breached the assumption that the non-zero effects come from a

303 Gaussian distribution. By design, the first three moments of the direct effects are fixed: they  
304 have zero mean, their variance is defined by the direct heritabilities and they must have zero  
305 skewness because the effect size distribution has to be symmetrical. Therefore, to violate the  
306 normality assumption, we varied the kurtosis (2, 3, 5, and 10) of the distribution drawn from  
307 the Pearson's distribution family. Third, we tested the assumption of the direct effects on our  
308 traits coming from a two-component Gaussian mixture by introducing a third component and  
309 observing how the estimates were effected. In this simulation scenario we introduced a large  
310 effect third component for  $X$  while decreasing the polygenicity of  $U$  ( $\pi_{x1} = 1 \times 10^{-4}$ ,  $\pi_{x2} =$   
311  $1 \times 10^{-2}$ ,  $h_{x1}^2 = 0.15$ ,  $h_{x2}^2 = 0.1$ ,  $\pi_u = 1 \times 10^{-2}$ ).

## 312 2.9 Application to real summary statistics

313 Once we demonstrated favourable performance of our method on simulated data, we went on  
314 to apply LHC-MR to summary statics obtained from the UK Biobank and other meta-analytic  
315 studies (Table S1) in order to estimate pairwise bi-directional causal effect between 13 complex  
316 traits. The traits varied between conventional risk factors (such as low education, high body  
317 mass index (BMI), dislipidemia) and diseases (including diabetes and coronary artery disease  
318 among others). SNPs with imputation quality greater than 0.99, and minor allele frequency  
319 (MAF) greater than 0.5% were selected. Moreover, SNPs found within the human leukocyte  
320 antigen (HLA) region on chromosome 6 were removed due to the abundance of SNPs associated  
321 with autoimmune and infectious diseases as well as the complicated LD structure present in that  
322 region. For traits with total heritability below 2.5%, the outgoing causal effect estimates were  
323 ignored since instrumenting such barely heritable traits is questionable.

324 In order to perform LHC-MR between trait pairs, a set of overlapping SNPs was used as input  
325 for each pair. The effects of these overlapping SNPs were then aligned to the same effect allele  
326 in both traits. To decrease computation time further (while only minimally reducing power), we  
327 selected every 10th QC-filtered SNP as input for the analysis. We calculated regression weights  
328 using the UK10K panel, which may be sub-optimal for summary statistics not coming from the  
329 UK Biobank, but we have previously shown<sup>21</sup> that estimating LD in a ten-times larger data set  
330 (UK10K) outweighs the benefit of using smaller, but possibly better-matched European panel  
331 (1000 Genomes<sup>22</sup>).

332 We also ran IVW for each trait pair in both directions to estimate bi-directional causal effects  
333 as well as LD score regression to get the cross trait intercept term. We then added uniformly  
334 distributed ( $\sim U(-0.1, 0.1)$ ) noise to these pre-estimated parameters to generate starting points  
335 for the second step of the likelihood optimisation. These closer-to-target starting points did not  
336 change the optimisation results, simply sped up the likelihood maximisation and increased the  
337 chances to converge to the same (primary) optimum. The LHC-MR procedure was run for each  
338 pair of traits 100 times, each using a different set of randomly generated starting points within  
339 the ranges of their respective parameters. For the optimisation of the likelihood function (Eq.  
340 2), we used the R function 'optim' from the 'stats' R package<sup>23</sup>. Once we fitted this *complete*  
341 model estimating 11 parameters in two steps  $\{i_x, i_y, \pi_x, \pi_y, h_x^2, h_y^2, t_x, t_y, \alpha_{x \rightarrow y}, \alpha_{y \rightarrow x}, i_{xy}\}$ , we  
342 then ran block jackknife to obtain the SE of the parameters estimated in the second step:  
343  $\{h_x^2, h_y^2, t_x, t_y, \alpha_{x \rightarrow y}, \alpha_{y \rightarrow x}, i_{xy}\}$ .

344 To support the existence of the confounders identified by LHC-MR, we used EpiGraphDB<sup>24</sup><sup>25</sup>  
345 to systematically identify those potential confounders. The database provided for each potential  
346 confounder of a causal relationship, a causal effect on trait  $X$  and  $Y$  ( $r1$ , and  $r3$  in their  
347 notation), the sign of the ratio of which ( $sign(r3/r1)$ ) was compared to the sign of the LHC-MR  
348 estimated  $t_y/t_x$  values representing the strength of the confounder acting on the two traits. We  
349 restricted our comparison to the sign only, since the  $r1, r3$  values reported in EpiGraphDB are

350 not necessarily on the same scale.

## 351 2.10 Comparison against conventional MR methods and CAUSE

352 We compared the causal parameter estimates of the LHC-MR method to those of five conven-  
353 tional MR approaches (MR-Egger, weighted median, IVW, mode MR, and weighted mode MR)  
354 using a Z-test [26]. The 'TwoSampleMR' R package [27] was used to get the causal estimates for  
355 all the pairwise traits as well as their standard errors from the above-mentioned MR methods.  
356 The same set of genome-wide SNPs that were used by LHC-MR, were used as input for the  
357 package. SNPs associated with the exposure were selected to various degrees (for simulation we  
358 selected SNPs over a range of thresholds: absolute P-value  $< 5 \times 10^{-4}$  to  $< 5 \times 10^{-8}$ ), and SNPs  
359 more strongly associated with the outcome than with the exposure (P-value  $< 0.05$  in one-sided  
360 t-test) were removed. The default package settings for the clumping of SNPs ( $r^2 = 0.001$ ) were  
361 used and the analysis was run with no further changes. We tested the agreement between the  
362 significance and direction of our estimates and that of standard MR methods, with the focus  
363 being on finding differences in statistical conclusions regarding causal effect sizes.

364 We compared our causal estimates from all our simulation settings to the causal estimates  
365 obtained by running MR-RAPS [1] also using the 'TwoSampleMR' R package, once by using the  
366 entire set of SNPs, and another by filtering for SNPs with a significance threshold of  $< 5 \times 10^{-4}$ .  
367 We also compared both our simulation as well as real data results against those of CAUSE [10].  
368 We first generated simulated data under the LHC model and used them as input to estimate the  
369 causal effect using CAUSE. We then generated simulated data using the CAUSE framework and  
370 inputted them to LHC-MR (as well as standard MR methods) to estimate the causal parameters.  
371 Lastly, we compared causal estimates obtained for the 78 trait pairs (156 bi-directional causal  
372 effects) from LHC-MR to those obtained when running CAUSE.

## 373 3 Results

### 374 3.1 Overview of the method

375 We fitted an 11-parameter structural equation model (SEM) (Figure 1) to genome-wide sum-  
376 mary statistics of two studied complex traits in order to estimate bi-directional causal effects  
377 between them (for details see Methods). Additional model parameters represent direct heri-  
378 tabilities for  $X$  and  $Y$ , confounder effects, cross-trait and individual trait LD score intercepts  
379 and the polygenicity for  $X$  and  $Y$ . All SNPs associated with the heritable confounder ( $U$ ) are  
380 indirectly associated with  $X$  and  $Y$  with effects that are proportional (ratio  $t_y/t_x$ ). SNPs that  
381 are directly associated with  $X$  (and not with  $U$ ) are also associated with  $Y$  with proportional  
382 effects (ratio  $1/\alpha_{x \rightarrow y}$ ). Finally, SNPs that are directly  $Y$ -associated are also  $X$ -associated with  
383 a proportionality ratio of  $1/\alpha_{y \rightarrow x}$ . These three groups of SNPs are illustrated on the  $\beta_x$ -vs- $\beta_y$   
384 scatter plot (Figure S2). In simple terms, the aim of our method is to identify the different  
385 clusters, estimate the slopes and distinguish which corresponds to the causal- and confounder  
386 effects. In this paper, we focus on the properties of the maximum likelihood estimates (and their  
387 variances) for the bi-directional causal effects arising from our SEM.

### 388 3.2 Simulation results

389 We started off with a realistic simulation setting of 234,000 SNPs on chromosome 10 (LD patterns  
390 used from the UK10K panel) and 50,000 samples for both traits. Traits  $X$ ,  $Y$  and confounder  $U$   
391 had average polygenicity ( $\pi_x = 5 \times 10^{-3}$ ,  $\pi_y = 1 \times 10^{-2}$ ,  $\pi_u = 5 \times 10^{-2}$ ), with substantial direct  
392 heritability for  $X$  and  $Y$  ( $h_x^2 = 0.25$ ,  $h_y^2 = 0.2$ ), mild confounding ( $t_x = 0.16$ ,  $t_y = 0.11$ ) and a

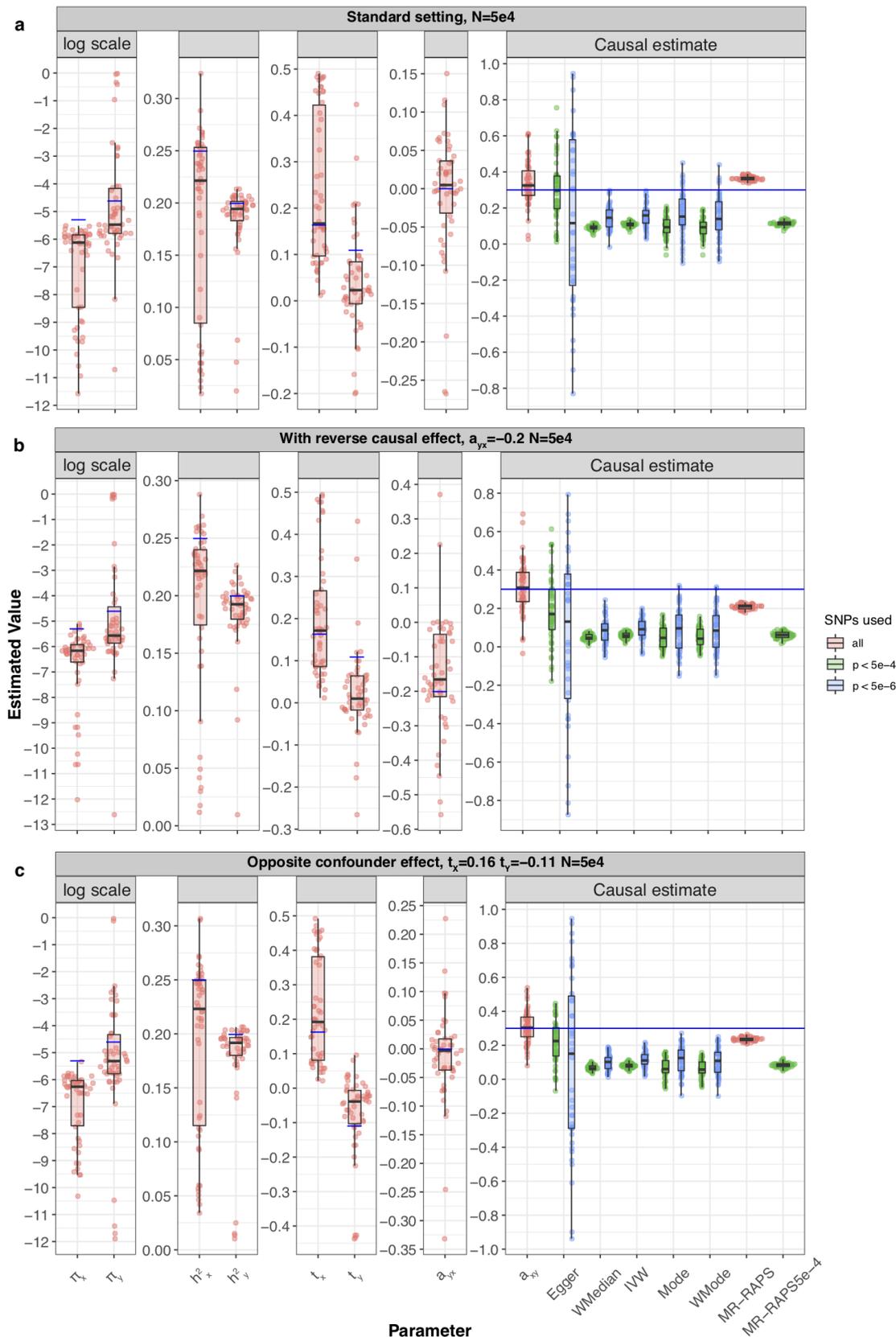
causal effect between  $X$  and  $Y$  ( $\alpha_{x \rightarrow y} = 0.3, \alpha_{y \rightarrow x} = 0$ ). Note that with these settings, SNPs associated with  $U$  would violate the InSIDE assumption but might still be used by conventional MR methods. Under this standard setting, there were no genome-wide significant SNPs for standard MR methods, and estimates derived using SNPs with a p-value  $< 5 \times 10^{-6}$  showed a downward bias for all MR methods (Figure 2 panel *a*). MR-RAPS using filtered SNPs (p-value  $< 5 \times 10^{-4}$ ) was similarly downward biased whereas MR-RAPS using the entire set of SNPs was upward biased with the least amount of variance compared to all methods including LHC-MR. LHC-MR in this scenario slightly over estimated the causal effect in comparison but had the smallest RMSE after MR-RAPS (0.13 vs 0.06, Supplementary Table S2).

We ran all our simulation scenarios with a smaller and a larger sample size (50,000 and 500,000) and observed that the relative performance of the methods were in some cases sample size specific. Smaller sample sizes often meant that standard MR methods had little to no IVs reaching genome-wide (GW) significance and hence we were forced to use IVs from less stringent thresholds ( $< 5 \times 10^{-4}$  and  $< 5 \times 10^{-6}$ ). Therefore, the causal effects were estimated with a substantial downward bias due to weak instrument bias (and winner's curse). LHC-MR in these cases was able to estimate the causal effect with less bias but with a larger variance compared to most standard MR methods – still outperforming them in terms of RMSE in most settings. In the larger sample size setting, standard MR methods had IVs for every threshold cutoff. However, a pattern also observed with smaller sample sizes, but to a lesser extent, the causal estimates of some methods changed (either in mean or in variance, most noticeably observed in weighted median and IVW) as the threshold became more stringent. This is of particular concern and highlights that while in this simulation setting the  $5 \times 10^{-8}$  threshold may have optimally cancelled out the different biases for IVW (downward bias due to winner's curse and weak instrument bias, upward bias due to genetic confounding), its estimate remains strongly setting-dependent. LHC-MR has performed reasonably well, exhibiting lower RMSE than most other methods, except for IVW and MR-RAPS for the  $5 \times 10^{-4}$  threshold (Figure S4 panel *a*). However, we observed that the performance of MR-RAPs is particularly setting- and threshold dependent.

Furthermore, unequal sample sizes for the two traits showed an underestimation of the causal effects for almost all MR methods, while LHC-MR remained the most accurate in the case where  $n_x$  (50,000) was smaller than  $n_y$  (500,000). However, the performances in the reverse scenario, where  $n_x$  was larger in size, were akin to the large sample size standard setting, where only IVW and filtered MR-RAPS ( $< 5 \times 10^{-4}$ ) showed superior performance to LHC-MR both in terms of bias and variance (see Figure S5).

When testing scenarios in the absence of a causal- or a confounder effect (imitating the classical MR assumptions), with a smaller causal effect ( $\alpha_{x \rightarrow y} = 0.1$ ), or with both forward- and reverse causal effects, we note that LHC-MR outperforms the standard MR methods as well as MR-RAPS in all these scenarios.

When there was no causal effect ( $\alpha_{x \rightarrow y} = 0$ ), LHC-MR had the smallest bias out of all the methods in both sample sizes (0.004 in both, Figure S6 panel *a* and Figure S7 panel *a*). The variance of the LHC-MR estimates in the larger sample size was much lower (0.0001 vs 0.01), similarly the other methods had a smaller variance in larger sample sizes and had more clearly seen upward biased estimates. The increased upward bias of standard MR methods is due to the fact that confounder-associated SNPs could only be detected in larger sample size and those lead to positive bias (due to the concordant effect of the confounder on the two traits). Note that the variance of standard MR methods are low simply because, in these settings, we were forced to lower the instrument selection threshold, hence artificially included many (potentially invalid) instruments, which lowers the estimator variance while increasing bias. MR-RAPS greatly overestimates the causal effects when the sample size is larger.



**Figure 2: Simulation results under various scenarios.** These Raincloud boxplots<sup>28</sup> represent the distribution of parameter estimates from 50 different data generations under various conditions. For each generation, standard MR methods as well as our LHC-MR were used to estimate a causal effect. The true values of the parameters used in the data generations are represented by the blue dots/lines. **a** Estimation under standard settings ( $\pi_x = 5 \times 10^{-3}$ ,  $\pi_y = 1 \times 10^{-2}$ ,  $\pi_u = 5 \times 10^{-2}$ ,  $h_x^2 = 0.25$ ,  $h_y^2 = 0.2$ ,  $h_u^2 = 0.3$ ,  $t_x = 0.16$ ,  $t_y = 0.11$ ). **b** Addition of a reverse causal effect  $\alpha_{y \rightarrow x} = -0.2$ . **c** Confounder with opposite causal effects on  $X$  and  $Y$  ( $t_x = 0.16$ ,  $t_y = -0.11$ ).

442 In the absence of a confounder effect, there is not much of a difference between the two sample  
443 sizes; standard MR methods have a large variance and are downward biased, LHC-MR is less  
444 biased compared to them but MR-RAPS performs best with the least bias and variance when all  
445 the SNPs are used as instruments (Figure S6 panel *b* and Figure S7 panel *b*). Trying a smaller  
446 causal effect led to an upward bias for all MR methods including both filterings of MR-RAPS  
447 in the larger sample size. Alternately, when  $n_x = n_y = 50,000$ , the MR methods are downward  
448 biased (Figures S6 panel *c* and Figures S7 panel *c*). Lastly, when a (negative) reverse causal  
449 effect is introduced, all MR methods and MR-RAPS are negatively biased in their estimation  
450 of the causal effect (see Figure 2 panel *b*). LHC-MR has a much smaller bias for the forward  
451 causal effect estimate in this case, and a generally small bias for the reverse causal effect in both  
452 sample sizes (0.05 for  $n = 50,000$  and 0.03 for  $n = 500,000$ , Figure S4 panel *b*).

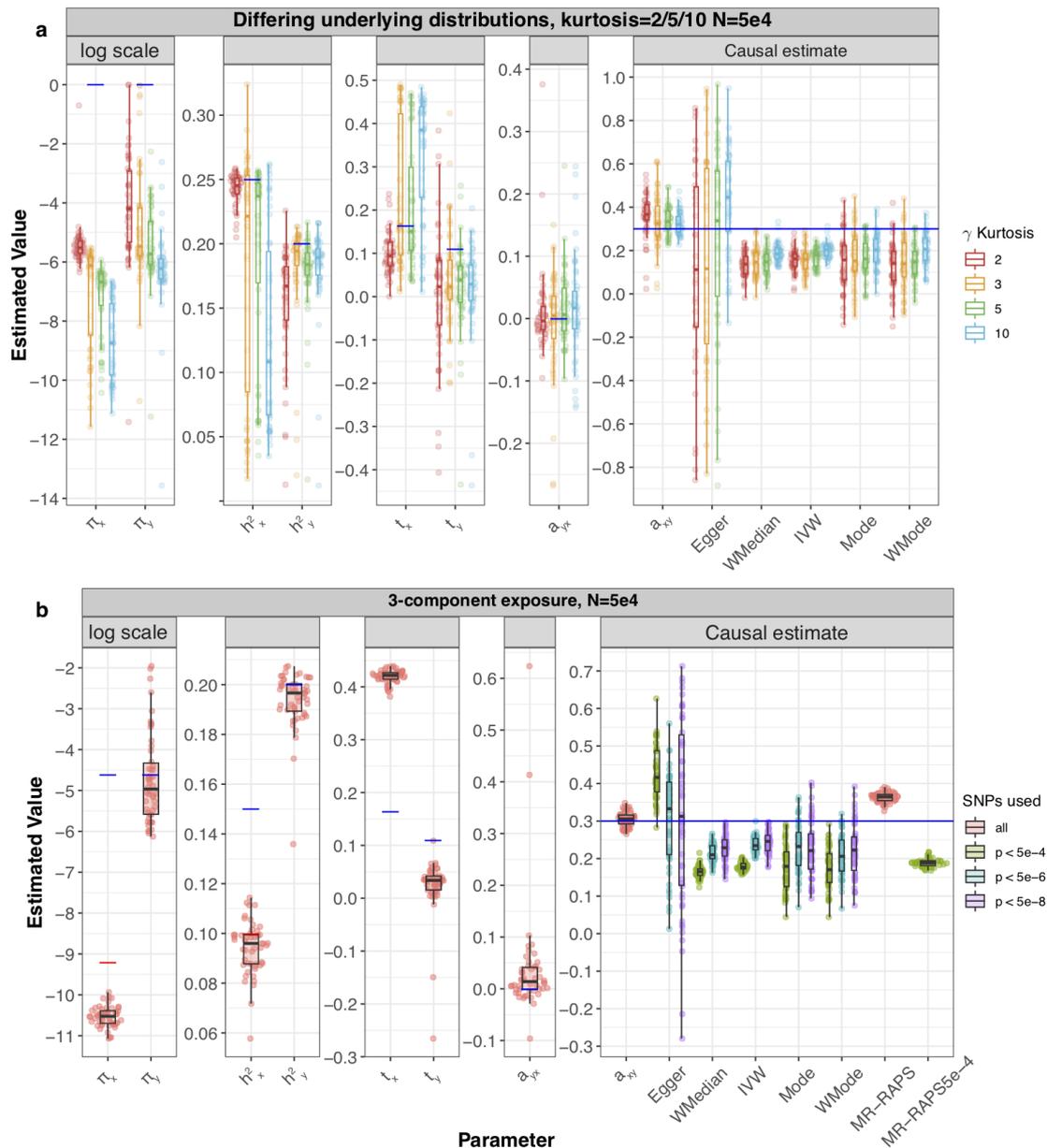
453 Increasing the indirect genetic effects, by intensifying the contribution of the confounder to  $X$   
454 and  $Y$  ( $t_x = 0.41, t_y = 0.27$ ), led to a general over estimation of the causal effects by all methods  
455 including LHC-MR, though more drastically seen in standard MR methods and MR-RAPS in  
456 larger sample sizes, when there is sufficient power to pick up these confounder-associated SNPs.  
457 The causal effect estimates of standard MR methods in the smaller sample size were much less  
458 affected by the presence of a strong confounder compared to LHC-MR and MR-RAPS (Figure  
459 S8). The reason for this is that the confounder-associated SNPs remain undetectable at lower  
460 sample size and hence instruments will not violate the classical MR assumptions.

461 Further testing the effects of the confounder trait on the causal estimation, we tested the impact  
462 of confounders with opposite effects on  $X$  and  $Y$ . We observe a major underestimation of the  
463 causal effects for standard MR methods as well as MR-RAPS, whereas LHC-MR performs better  
464 for both sample sizes (RMSE = 0.01 and 0.1 for larger and smaller  $n$  respectively), see Figures  
465 2 panel *c* and S4 panel *c*.

466 Our LHC-MR method is influenced by the unlikely scenario of extreme polygenicity for traits  
467  $X, Y$  and  $U$ , and it suffers from increased bias and variance regardless of sample size (see Figure  
468 S9). Standard MR methods as well as filtered MR-RAPS underestimated the causal effect  
469 when  $n = 50,000$ . Some also underestimated  $\alpha_{x \rightarrow y}$  when  $n = 500,000$ , with the exception of  
470 IVW, Mode and filtered MR-RAPS, that outperformed the rest. Decreasing the proportion of  
471 confounder-associated SNPs to 1% only, does not seem to affect our method and shows similar  
472 results to the standard setting (Figure S10).

473 Furthermore, we simulated summary statistics, where (contrary to our modelling assumptions)  
474 the  $X - Y$  relationship has two confounders,  $U_1$  and  $U_2$ . When the ratio of the causal effects  
475 of these two confounders on  $X$  and  $Y$  ( $q_x^{(1)}/q_y^{(1)}$  and  $q_x^{(2)}/q_y^{(2)}$  respectively) agreed in sign, the  
476 corresponding causal effects of standard MR methods were over-estimated in larger sample sizes  
477 and, conversely, underestimated in smaller sample sizes (Figures S11 and S12, panels *a*). LHC-  
478 MR and weighted median performed better however in larger sample sizes and had a bias of 0.03  
479 and 0.07 respectively. However, when the signs were opposite ( $q_x^{(1)} = 0.3, q_y^{(1)} = 0.2$  for  $U_1$  and  
480  $q_x^{(2)} = 0.3, q_y^{(2)} = -0.2$  for  $U_2$ ), conventional MR methods and MR-RAPS in this case almost all  
481 underestimated the causal effect regardless of sample size. LHC-MR outperformed them both  
482 in the larger sample size (bias of 0.007) and in the smaller sample size (bias of  $-0.003$ ), see  
483 Figures S11 and S12, panels *b*.

484 Finally, we explored how sensitive our method is to different violations of our modelling assump-  
485 tions. First, we simulated summary statistics when the underlying non-zero effects come from a  
486 non-Gaussian distribution. Interestingly, we observed that, for smaller sample sizes, the variance  
487 of the causal effect estimate was dependent on the kurtosis for most MR methods. LHC-MR  
488 estimations yielded slightly more pronounced upward bias than IVW, while still exhibiting the  
489 lowest RMSE among all methods (Figure 3 panel *a*). Similar results are seen in larger sample



**Figure 3: Simulation results under various scenarios.** These Raincloud boxplots<sup>28</sup> represent the distribution of parameter estimates from 50 different data generations under various conditions. For each generation, standard MR methods as well as our LHC-MR were used to estimate a causal effect. The true values of the parameters used in the data generations are represented by the blue dots/lines. **a** The different coloured boxplots represent the underlying non-normal distribution used in the simulation of the three  $\gamma_x, \gamma_x, \gamma_u$  vectors associated to their respective traits. The Pearson distributions had the same zero mean and skewness, however their kurtosis ranged between 2 and 10, including the kurtosis of 3, which corresponds to a normal distribution assumed by our model. The standard MR results reported had IVs selected with a p-value threshold of  $5 \times 10^{-6}$ . **b** Addition of a third component for exposure  $X$ , while decreasing the strength of  $U$ . True parameter values are in colour, blue and red for each component ( $\pi_{x1} = 1 \times 10^{-4}, \pi_{x2} = 1 \times 10^{-2}, h_{x1}^2 = 0.15, h_{x2}^2 = 0.1$ ).

size with smaller variance for all methods under all degrees of kurtosis except for IVW, which showed a better performance than LHC-MR (Figure S13 panel *a*). Second, we simulated effect sizes coming from a three-component Gaussian mixture distribution (null/small/large effects), instead of the classical spike-and-slab assumption of our model. The smaller sample size estimates mirror those of the standard setting with  $n$  also equal to 50,000 (see Figure 3 panel *b*). However, in the larger sample size, LHC-MR overestimates the causal effect. This bias could be due to the merging of true effect estimates with confounder effect leading to an overestimation of  $\alpha_{x \rightarrow y}$  (Figure S13 panel *b*). MR-Egger, IVW and filtered MR-RAPS have the smallest RMSE in this case.

### 3.2.1 Comparing CAUSE and LHC-MR

When running CAUSE on data simulated using the LHC-MR model framework in order to estimate a causal effect ( $\gamma$  in their notation), we investigated three different scenarios, each with multiple data generations: one where the underlying model has a shared factor/confounder with effect on both exposure and outcome only, another where the underlying model has a causal effect of 0.3 only, a third where the underlying model has both a causal effect and a shared factor. The data generated using the LHC-MR model was done under the standard settings ( $\pi_x = 5 \times 10^{-3}$ ,  $\pi_y = 1 \times 10^{-2}$ ,  $\pi_u = 5 \times 10^{-2}$ ,  $h_x^2 = 0.25$ ,  $h_y^2 = 0.2$ ,  $h_u^2 = 0.3$ ,  $t_x = 0.16$ ,  $t_y = 0.11$ ,  $\alpha_{x \rightarrow y} = 0.3$ ,  $\alpha_{y \rightarrow x} = 0$ ,  $m = 234,000$ ,  $n_x = n_y = 50,000$ ). For each setting, 50 different replications were investigated.

In the case of an underlying shared effect only, CAUSE preferred the sharing model 100% of the time, and thus there was no causal estimation, however it underestimated both  $\eta$  and  $q$ . When there was an underlying causal effect only, CAUSE preferred the causal model only 4% of the times, where it slightly underestimated the causal effect ( $\hat{\gamma} = 0.241$ ). Although the true values of  $\eta$  and  $q$  are null in this scenario, the sharing model returned estimates for these two parameters overestimating them both (probably driven by their priors), as seen in Figure S14. In the third case, and in the presence of both, CAUSE preferred the sharing model in 48 of the 50 simulations, yet it underestimated  $\eta$  (corresponding to  $t_y/t_x$  for our model) but overestimated  $q$  ( $t_x^2/(t_x^2 + h_x^2)$  in our model) (mean of 0.566 and 0.222 respectively where the true values are 0.667 and 0.097) showing a similar estimation pattern to the second case. Interestingly, in larger sample sizes, CAUSE selects the correct model 100% of the time, but still underestimates  $\gamma$ , see Figure S15.

In the reverse situation, where data was generated using the CAUSE framework (with parameters  $h_1 = h_2 = 0.25$ ,  $m = 97,450$ ,  $N1 = N2 = 50,000$ ) and LHC-MR was used to estimate the causal effect, we saw the following results (see Figure S16). First, when we generated data in the absence of causal effect ( $\gamma = 0$ ,  $\eta = \sqrt{0.05}$ ,  $q = 0.1$ ), CAUSE does extremely well in estimating a null causal effect 100% of the time. Standard MR methods yield a slight overestimation of the (null) causal effect with varying degrees of variance, whereas LHC-MR shows both a greater variance and an upward bias – still leading to a causal effect compatible with zero. Second, in the absence of a confounder combined with non-zero causal effect ( $\gamma = \sqrt{0.05} = 0.22$ ,  $\eta = 0$ ,  $q = 0$ ), CAUSE underestimates the causal effect ( $\hat{\gamma} = 0.18$ ) compared to LHC-MR which overestimates the causal effect: the mean of the estimates was 0.38 (over the 50 runs). Finally, in the presence of both a confounder and a causal effect ( $\gamma = \sqrt{0.05}$ ,  $\eta = \sqrt{0.05}$ ,  $q = 0.1$ ), CAUSE slightly underestimates the causal effect ( $\hat{\gamma} = 0.20$ ), whereas LHC-MR overestimates the effects and shows estimates reaching the boundaries 11 out of 50 times (mean of the converged  $\hat{\gamma} = 0.39$  over the 39 data simulations, see Figure S16 panel *c*) – indicating that this setting of the CAUSE model is not compatible with the LHC-MR model framework. Interestingly, classical MR methods outperform CAUSE in this case. Note that in the interest of run time we used less

537 SNPs (than usual) for parameter estimations. The analysis was repeated for a larger sample size  
538 of 500,000 (Figure S17), with more favourable results for LHC-MR. In the absence of a causal  
539 effect, we had similar results to smaller sample sizes, whereas in the absence of a shared effect,  
540 LHC-MR estimates the causal effect accurately with a mean of 0.22, CAUSE underestimates  
541 it and the rest of the MR methods are less biased. In the presence of both causal and shared  
542 factor, CAUSE recovers the causal effect. IVW, unlike the other MR methods and CAUSE, is  
543 more affected by the presence of the confounder, while LHC-MR exhibits upward bias with a  
544 mean estimate of 0.27.

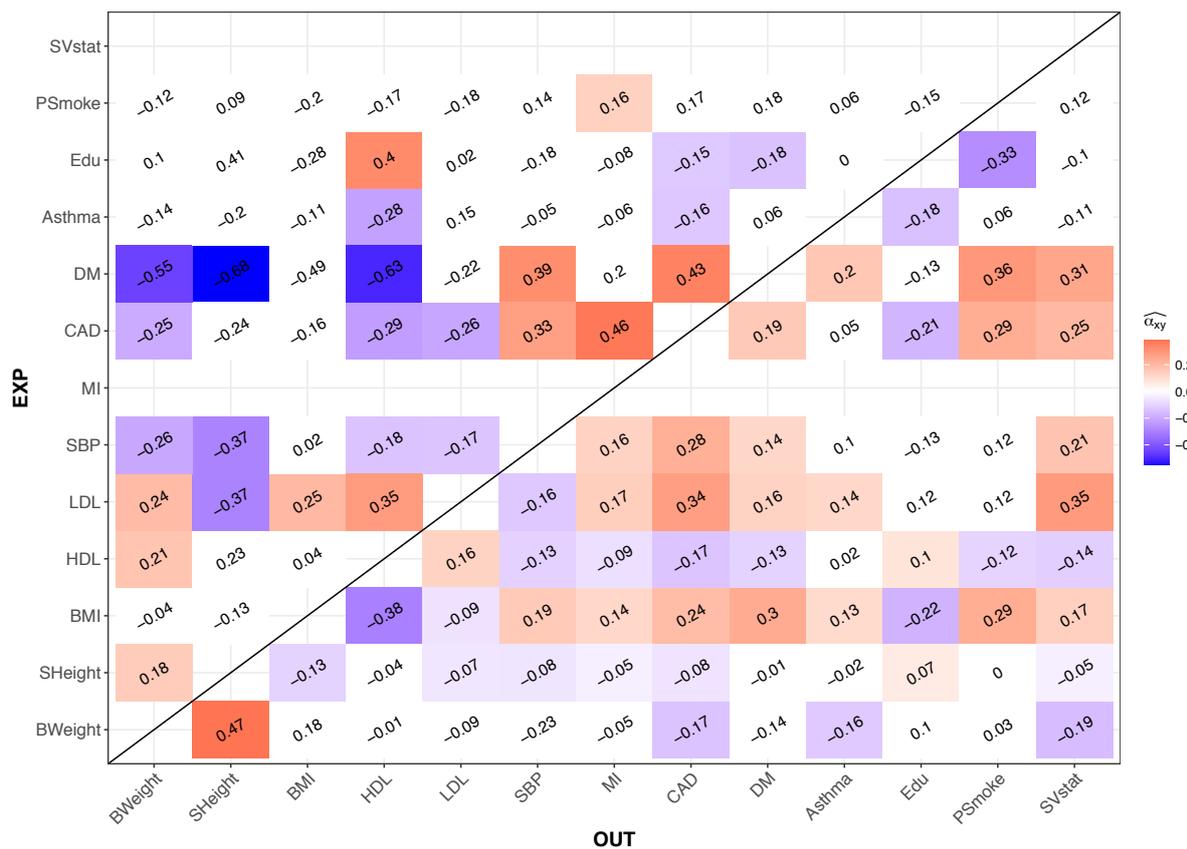
### 545 3.3 Application to association summary statistics of complex traits

546 We applied our LHC-MR and other MR methods to estimate all pairwise causal effects between  
547 13 complex traits (156 causal relationships in both directions). Our results are presented as a  
548 heatmap in Figure 4 (and are detailed in Supplementary Table S3). Further, we calculated the  
549 alternate set of estimated parameters that naturally results from our model (for reference see  
550 Sections 2.2 and Supplementary materials 1.1). Among trait pairs for which the exposure had  
551 sufficient heritability ( $> 2.5\%$ ), the alternate parameters of a 102 trait pairs were within the  
552 possible ranges mentioned in methods (i.e. the confounder and the exposure are interchange-  
553 able). However, for all of these pairs, the alternative parameter optima lead to lower direct-  
554 than indirect heritability, which we deem unrealistic. Therefore, we report only the primary set  
555 of estimated optimal parameters in the main results and provide the alternative parameters in  
556 the Supplementary Table S4. The comparison of the results obtained by LHC-MR and stan-  
557 dard MR methods is detailed below and more extensively in Supplementary Tables S5-S6. In  
558 summary, LHC-MR provided reliable causal effect estimates for 132 out of 156 exposure traits  
559 (i.e. those exposures had an estimated total heritability greater than 2.5%). These estimates  
560 were compared to five different MR methods. Seventy-four causal relationships were deemed  
561 significant by LHC-MR. Furthermore, for 117 out of those 132 comparable causal relationships,  
562 our LHC-MR causal effect estimates were concordant (not significantly different) with at least  
563 two out of five standard MR methods' estimates.

564 By simply comparing the significance status and the direction of the causal effects between the  
565 methods, we see that LHC-MR agrees in sign and significance (or the lack there of) with at least  
566 3 MR methods 77 times. For 31 relationships, LHC-MR results lead to different conclusions  
567 than those of standard MR methods. For 28 of those, LHC-MR identified a causal effect missed  
568 by all standard MR methods. For the other three, we observed a disagreement in sign: LDL  
569 has a negative effect on BMI according to weighted mode and weighted median, whereas we  
570 show a positive effect, HDL and LDL show a negative bi-directional causal effect for weighted  
571 mode but a positive bi-directional effect with LHC-MR. Despite the conflicting evidence for the  
572 causal relationship of LDL on BMI, studies have shown that the relationship between them is  
573 non-linear<sup>29</sup>, possibly explaining the discrepancy between the results.

574 LHC-MR agreed with most MR estimates and confirmed many previous findings, such as in-  
575 creased BMI leading to elevated blood pressure<sup>30,31</sup>, diabetes mellitus<sup>32,33</sup> (DM), myocardial  
576 infarction<sup>34</sup> (MI) and coronary artery disease<sup>35</sup> (CAD). Furthermore, we confirmed previous  
577 results<sup>36</sup> that diabetes increases SBP ( $\hat{\alpha}_{x \rightarrow y} = 0.39 - P = 1.70 \times 10^{-9}$ ).

578 Interestingly, it revealed that higher BMI increases smoking intensity, concordant with other  
579 studies<sup>37,38</sup>. It has also shown the protective effect of education against a range of diseases  
580 (e.g. CAD and diabetes<sup>39,40</sup>) and risk factors such as smoking<sup>41,42</sup>, in agreement with previous  
581 observational and MR studies. Probably reflecting lifestyle change recommendations by medical  
582 doctors upon disease diagnosis, statin use is greatly increased when being diagnosed with CAD,  
583 (systolic) hypertension, dislipidemia, and diabetes as is shown by both LHC-MR and standard  
584 MR methods.



**Figure 4: Heatmap representing the bi-directional causal relationship between the 13 UK Biobank traits.** The causal effect estimates in coloured tiles all have a significant p-values surviving Bonferroni multiple testing correction with a threshold of  $3.2 \times 10^{-4}$ . We did not report an estimated causal effects for exposures with an estimated total heritability less than 2.5%. White tiles show an absence of a significant causal effect estimate.

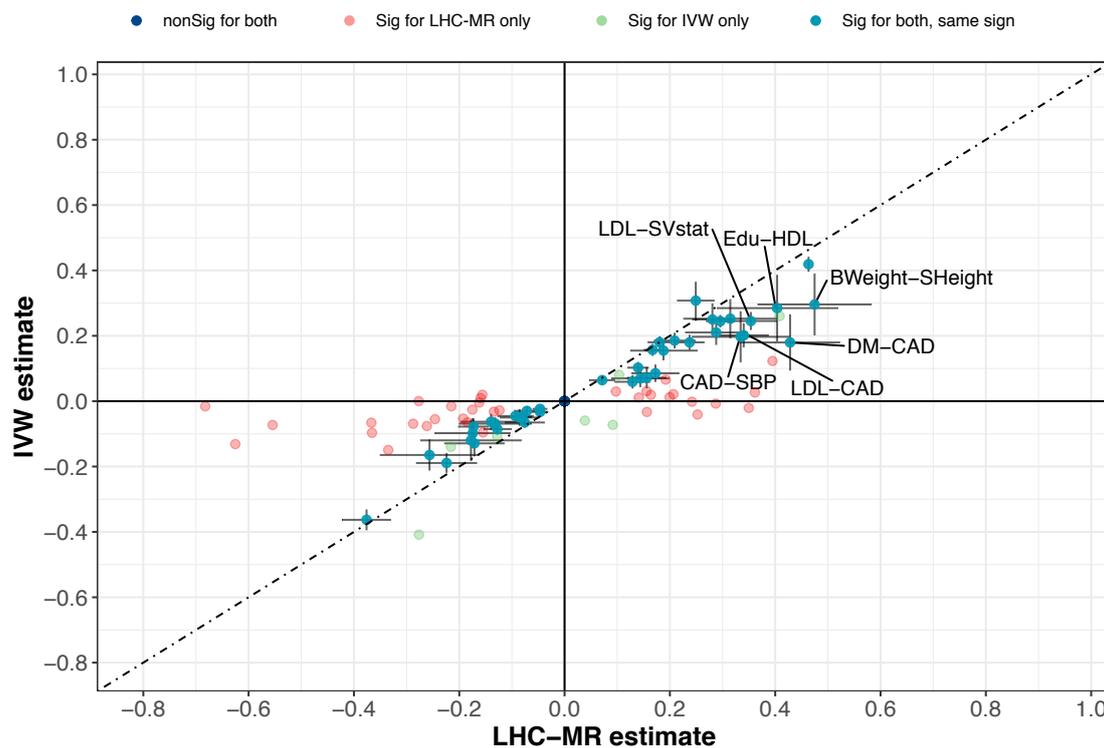
Abbreviations: BMI: Body Mass Index, BWeight: Birth Weight, CAD: coronary artery disease, DM: Diabetes Mellitus, Edu: Years of Education, HDL: High-density Lipoprotein LDL: Low-density Lipoprotein, MI: Myocardial Infarction, PSmoke: # of Cigarettes Previously Smoked, SBP: Systolic Blood Pressure, SHHeight: Standing Height, SVstat: Medication Simvastatin

585 Furthermore, causal effects of height on CAD, DM and SBP have been previously examined in  
 586 large MR studies [43,44]. LHC-MR, agreeing with these claims, did not find significant evidence to  
 587 support the effect of height on DM, but did find a significant protective effect on CAD and SBP.  
 588 However, unlike the first two, the relationship between height and SBP also revealed the existence  
 589 of a confounder with causal effects 0.14 ( $P = 9.2 \times 10^{-11}$ ) and 0.11 ( $P = 3.39 \times 10^{-8}$ ) on height  
 590 and SBP respectively. Another example of a trait pair for which LHC-MR found an opposite  
 591 sign confounder effect is HDL and its protective effect on SBP. The confounder had a positive  
 592 effect ratio of  $t_y/t_x = 0.84$ , opposing the negative causal effect of  $\hat{\alpha}_{x \rightarrow y} = -0.13$  supported by  
 593 observational studies [45]. This causal effect was not found by any other MR method.

594 It is important to note that while the effects of parental exposures on offspring outcomes can be  
 595 seen as genetic confounding, LHC-MR would not be able to distinguish parental and offspring  
 596 causal effects, because the LHC-MR model assumes that there is no correlation between the  
 597 genetic effects on the exposure and the genetic effects on the confounder (which is not the  
 598 case of parental vs offspring traits). Thus, LHC-MR causal effect estimates are just as likely  
 599 to reflect parental effects as any other MR method [46]. This may be the case, for example,

600 for the detrimental effect of increased (parental) BMI on education (supported by longitudinal  
601 studies<sup>[47]</sup>), the positive effect of (parental) height on birth weight<sup>[48]</sup>, or on education<sup>[49]</sup>. There  
602 are also some associations identified only by LHC-MR that might reflect parental effects: the  
603 negative causal effect of CAD on education or on birth weight, the positive impact of HDL on  
604 birth weight, or DM reducing height. All these pair associations uniquely found by LHC-MR  
605 are examples of LHC-MR's use of whole-genome SNPs instead of GW-significant SNPs only,  
606 as our estimates are of larger magnitude than those found by standard MR. Interestingly, for  
607 the CAD→birth weight relationship, LHC-MR revealed a confounder of opposite causal effects,  
608 which could have masked/mitigated the causal effect of standard MR methods.

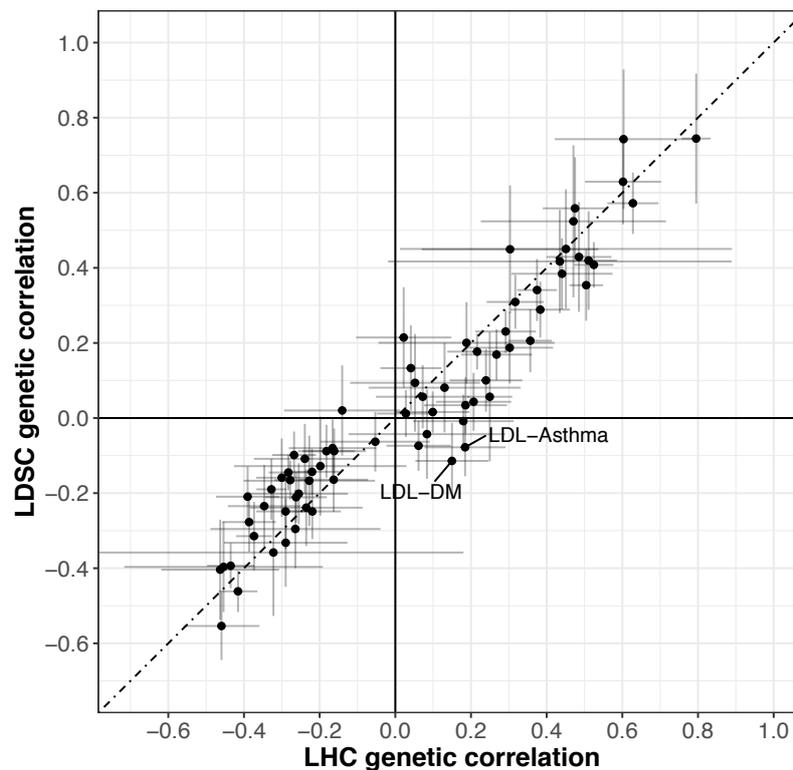
609 A systematic comparison between IVW and LHC-MR has shown generally good agreement  
610 between the two methods, which is illustrated in Figure 5. To identify discrepancies between  
611 our causal estimates and those of the standard MR results, we grouped the estimates into several  
612 categories, either non-significant p-value for both or either, significant with an agreeing sign for  
613 the causal estimate, or significant with a disagreeing sign. The diagonal (seen in Figure 5)  
614 representing the agreement in significance status and sign between the two methods, is heavily  
615 populated. On the other hand, 34 pairs have causal links that are significantly non-zero according  
616 to LHC-MR, but are non-significant for IVW, while the opposite is true for seven pairs. We  
617 believe that many of these seven pairs may be false positives, since four of them are picked up  
618 by no other MR method, two are confirmed by only one other method and the last one by two  
619 methods. Further comparisons of significance between LHC-MR estimates and the remaining  
620 standard MR methods can be found in Table S7.



**Figure 5:** A scatter plot of the causal effect estimates between LHC-MR and IVW. To improve visibility, non-significant estimates by both methods are placed at the origin, while significant estimates by both methods appear on the diagonal with 95% CI error bars.

621 LHC-MR identified a confounder for 16 trait pairs out of the possible 78. In order to support

622 these findings, we used EpiGraphDB<sup>[24,25]</sup> to systematically identify those potential confounders.  
623 EpiGraphDB could identify reliable confounders for ten out of the 16 trait pairs. Notably, for  
624 the birth weight - diabetes pair, the average epigraph confounder-effect ratio ( $r_3/r_1$ ) clearly  
625 agreed in sign with our  $t_y/t_x$  ratio, indicating that the characteristics of the confounder(s)  
626 evidenced by LHC-MR agree with those found in an exhaustive confounder search, and are  
627 mainly obesity-related traits (Figure S18 panel *a*). Six other trait pairs showed mixed signs of  
628 different confounders, indicating the possibility of having heterogeneous confounders (Figure S18  
629 panels *b-e*). Finally, three trait pairs showed a disagreement between our estimated confounder  
630 effect ratio and the bulk of those found by epigraphDB as seen in Supplementary Figure S18  
631 panels *f-j*. However, at least one of the top ten potential confounders showed effects that are in  
632 agreement with our ratio for each of these pairs. Note that since the reported causal effects of  
633 the confounders on  $X$  and  $Y$  reported in EpiGraphDB are not necessarily on the same scale, we  
634 do not expect the magnitudes to agree.



**Figure 6:** Scatter plot comparing the genetic correlation for each trait obtained from LDSC against the value calculated using parameter estimates from the LHC-MR model. A 95% CI is shown for each point. Values from both methods are reported in Supplementary Table S9.

635 As described in the methods (Eq. 3), genetic correlation can be computed from our estimated  
636 model parameters. To verify that the fitted LHC-MR model leads to a genetic correlation simi-  
637 lar to the one obtained from LD score regression<sup>[50]</sup> (LDSC), we compared whether the two  
638 approaches produce similar genetic correlation estimates. We did this by taking the estimated  
639 parameters obtained from the 200 block jackknife to estimate the genetic correlations between  
640 traits (and their standard errors), and plotted them against LD score regression values as seen  
641 in Figure 6. As expected, we observe an overall good agreement between the estimates of the  
642 two methods, with only six trait pairs differing in sign. Of these six, only 2 were nominally  
643 significantly different between the two methods (LDL→Asthma and LDL→DM). Further de-  
644 composition of the genetic covariance into heritable confounder-led or causal effect-led revealed

645 that most of the genetic covariance between traits can be attributed to bi-directional causal  
646 effects. A reason for this could be that confounders would need to have very strong effects  
647 to substantially contribute to the genetic correlation ( $\approx t_x \cdot t_y$ ) compared to the bi-directional  
648 causal effects ( $\approx \alpha_{x \rightarrow y}^2 \cdot h_x^2 + \alpha_{y \rightarrow x}^2 \cdot h_y^2$ ).

649 As for the comparison of LHC-MR against CAUSE for real trait pairs, we ran CAUSE on all  
650 156 trait pairs (bi-directional), and extracted the parameter estimates that corresponded to  
651 the methods winning model. The p-value threshold was corrected for multiple testing and was  
652 equivalent to 0.05/156. Based on that threshold, the p-value that compared between the causal  
653 and the sharing model of CAUSE was used to choose one of the two. Then the parameters  
654 estimated from the winning model,  $\gamma$  (only for causal model),  $\eta$  and  $q$ , were compared to their  
655 counterparts in LHC-MR. A visual comparison of LHC-MR's causal estimates and those of  
656 CAUSE can be seen in Figure [S19](#).

657 Whenever the causal effect estimates were significant both for CAUSE and LHC-MR (30 causal  
658 relationships), they always agreed in sign (Table [S8](#)) with a high Pearson correlation of 0.592.  
659 Calculating the correlation for their estimates regardless of significance yielded a smaller value of  
660 0.377. When compared to the causal effect estimate from IVW, LHC-MR was strongly correlated  
661 (0.585), whereas CAUSE had a slightly weaker correlation (0.471) using all estimates.

662 Similarly, the significant confounder effect ratio of LHC-MR ( $t_y/t_x$ ) can be compared to the  
663 significant confounder effect estimate of CAUSE ( $\eta$ ) when a sharing model is chosen. These  
664 12 confounding quantities by CAUSE and LHC-MR disagreed in sign for all but one trait pair  
665 (Height $\rightarrow$ MI), with a Pearson correlation compatible with zero ( $-0.357$  (95% CI  $[-0.77, 0.27]$ ))  
666 .

## 667 4 Discussion

668 We have developed a structural equation (mixed effect) model to account for a latent heritable  
669 confounder ( $U$ ) of an exposure ( $X$ ) - outcome ( $Y$ ) relationship in order to estimate bi-directional  
670 causal effects between the two traits ( $X$  and  $Y$ ). The method, termed LHC-MR, fits this model  
671 to association summary statistics of genome-wide genetic markers to estimate various global  
672 characteristics of these traits, including bi-directional causal effects, confounder effects, direct  
673 heritabilities, polygenicities, and population stratification.

674 We first demonstrated through simulations that in most scenarios, the method produces causal  
675 effect estimates with substantially less bias and variance (in larger sample sizes) than other MR  
676 tools. The direction and magnitude of the bias of classical MR approaches varied across scenarios  
677 and sample sizes. This bias was mainly influenced by two often opposite forces: downward  
678 bias resulting from winner's curse and weak instruments, and upward bias due to a positive  
679 confounder of the  $X - Y$  relationship, evident in larger sample sizes. In the scenario lacking a  
680 confounder (thus respecting all MR assumptions), MR methods were distinctly underestimating  
681 the causal effect, except for LHC-MR and to a better extent MR-RAPS. However, under standard  
682 settings with an added small heritable confounder and no reverse causality present, all classical  
683 MR methods still slightly underestimated the causal effect in smaller sample sizes, except for the  
684 MR-RAPS estimate which was now overestimated. For the same standard setting scenario but  
685 in a larger sample size where confounder effects were more detectable, IVW had an estimation  
686 that was close to the true causal value chosen ( $\alpha_{x \rightarrow y} = 0.3$ ) due to the opposite biases cancelling  
687 out. However, when the causal effect was set to be smaller ( $\alpha_{x \rightarrow y} = 0.1$ ), the estimates of IVW  
688 became biased. More substantial violations of classical MR assumptions, such as the presence  
689 of negative-effect confounder or a negative reverse causal effect, led to more substantial biases  
690 that impacted all methods (including MR-RAPS) except LHC-MR.

691 Interestingly, in smaller sample sizes, standard MR methods showed a slight decreasing trend in  
692 the variance of the causal effect estimate as the kurtosis of the underlying effect size distribution  
693 went up from 2 to 10. On the other hand, LHC-MR did not show a similar trend with growing  
694 kurtosis, and estimated the causal effect with a smaller bias. As confounder causal effects ( $q_x$ ,  
695  $q_y$ ) increased, classical MR methods (except weighted ones) were prone to produce overestimated  
696 causal effects with at least twice the bias than that of LHC-MR, especially in larger sample sizes  
697 where the confounder-associated SNPs make it to the set of GW-significant instruments for all  
698 methods. Furthermore, mode-based estimators were robust to the presence of two concordant  
699 confounders, yet their bias was still 10-fold higher than LHC-MR's, and they did not perform  
700 as well in the presence of discordant confounders. In summary, LHC-MR was robust to a wide  
701 range of violations of the classical MR assumptions and was less impacted than standard MR  
702 methods. Thus it outperformed all MR methods in virtually all tested scenarios, many of which  
703 violated even its own modelling assumptions.

704 We then applied our method to summary statistics of 13 complex traits from large studies,  
705 including the UK Biobank. We observed a general trend in our results that (in agreement with  
706 epidemiological studies) higher BMI and LDL are risk factors for most diseases such as diabetes  
707 and CAD. We also note the protective effect HDL has on these same diseases. Moreover, we  
708 observe many disease traits increasing the intake of lipid-lowering medication (simvastatin),  
709 reflecting the recommendation/treatment of medical personnel following the diagnosis.

710 LHC-MR can have discordant results compared to other MR methods for many possible reasons.  
711 The positive causal effect of smoking on MI, diabetes on asthma, the protective impact of higher  
712 birth weight on asthma, or higher education on smoking intensity, all of which were missed by  
713 standard MR could reflect the increased power of LHC-MR with its use of full-genome SNPs as  
714 opposed to genome-wide significant SNPs of classical MR approaches. Estimates from classical

715 MR methods could also be impacted by sample overlap between the exposure and outcome data-  
716 sets, whereas LHC-MR takes this into account. However, when using large sample sizes, the bias  
717 due to sample overlap is expected to be very small, and therefore not sufficient to explain any  
718 discrepancy in the results<sup>51</sup>. Another possible reason for the discrepancy between our findings  
719 and those of standard MR methods is the presence of a significant heritable confounder found  
720 by LHC-MR with opposite effect to the estimated causal effect between the pair. These two  
721 opposite forces lead to association summary statistics that may be compatible with reduced (or  
722 even null) causal effect when the confounder is ignored. Possible examples of this scenario can  
723 be observed for when (parental) traits, e.g. diabetes and CAD, act on birth weight. These pairs  
724 have a confounder of opposite effects, possibly related to (parental) obesity. Similarly, standard  
725 MR methods show little evidence for a causal effect of SBP on height, while our LHC-MR  
726 estimate is  $-0.37$  ( $P = 4.81 \times 10^{-8}$ ) which most probably reflects parental (maternal) effects as  
727 seen in previous studies<sup>52</sup><sup>53</sup>. The protective effect of HDL on SBP is another example where a  
728 confounder of opposite sign to that of the causal effect allows it to be uniquely found by LHC-  
729 MR. LHC-MR assumes no genetic correlation between the confounder and the direct effects on  
730 the exposure, which may be violated when the confounder is the same trait as the exposure,  
731 but in the parent. Such parental effects can mislead most MR methods<sup>54</sup>, including ours,  
732 and hence we may observe biased results for traits such as BMI $\rightarrow$ education and HDL $\rightarrow$ birth  
733 weight.

734 Sixteen trait pairs showed a strong confounder effect, in the form of significant  $t_x$  and  $t_y$  esti-  
735 mates. These pairs were investigated for the presence of confounders using EpiGraphDB, and 10  
736 of them returned possible confounders. The bulk of such pairs returned confounders with both  
737 agreeing and disagreeing effect directions on  $X$  and  $Y$ , making it difficult to pinpoint a group of  
738 concordant and dominant confounders. However, for the birth weight-DM pair, where LHC-MR  
739 identifies a negative reverse causal effect and a confounder with effects  $t_x = 0.10$  ( $P = 6.77 \times 10^{-8}$ )  
740 and  $t_y = 0.15$  ( $P = 3.13 \times 10^{-7}$ ) on birth weight and DM respectively, EpiGraphDB confirmed  
741 several confounders related to body fat distribution and weight that matched in sign with our  
742 estimated confounder effect (Figure S18 panel *a*). Note that EpiGraphDB causal estimates are  
743 not necessarily on the scale of SD outcome difference upon 1 SD exposure change scale, hence  
744 they are not directly comparable with the  $t_y/t_x$  ratio, but are rather indicative of the sign of  
745 the causal effect ratio of the confounder. Furthermore, if EpiGraphDB does not find a causal  
746 relationship between the trait pair in either directions, then it does not return any possible con-  
747 founders of the two, a reason why only 10 out of 16 confounder-associated trait pairs returned  
748 any hits.

749 Lastly, our comparison of the genetic correlations calculated from our estimated parameters  
750 against those calculated from LD score regression showed good concordance, confirming that  
751 the detailed genetic architecture proposed by our model is compatible with the observed genetic  
752 covariance. The major difference between the genetic correlation obtained by LD score regres-  
753 sion *vs* LHC-MR is that our model approximates all existing confounders by a single latent  
754 variable, which may be inaccurate when multiple ones exist with highly variable  $t_y/t_x$  ratios.  
755 Furthermore, LHC-MR decomposed the observed genetic correlation into confounder and bi-  
756 directional causality driven components, revealing that most genetic correlations are primarily  
757 driven by bi-directional causal effects. Note that we have much higher statistical power to de-  
758 tect situations when the confounder effects are of opposite sign compared to the causal effects,  
759 because opposing genetic components are more distinct.

760 To our knowledge only two recent papers use similar models and genome-wide summary statis-  
761 tics. The LCV approach<sup>55</sup> is a special case of our model, where the causal effects are not  
762 included in the model, but they estimate the confounder effect mixed with the causal effect to

763 estimate a quantity of genetic causality proportion (GCP). In agreement with others<sup>[10,56]</sup>, we  
764 would not interpret non-zero GCP as evidence for causal effect. Moreover, in other simulation  
765 settings, LCV has shown very low power to detect causal effects (by rejecting  $GCP=0$ ) (Fig S15  
766 in Howey et al.<sup>[57]</sup>). Another very recent approach, CAUSE<sup>[10]</sup>, proposes a structural equation  
767 mixed effect model similar to ours. However, there are several differences between LHC-MR  
768 and CAUSE: (a) we allow for bi-directional causal effects and model them simultaneously, while  
769 CAUSE is fitted twice for each direction of causal effect; (b) they first use an adaptive shrinkage  
770 method to integrate out the multivariable SNP effects and then go on to estimate other model  
771 parameters, while we fit all parameters at once; (c) CAUSE estimates the correlation parameter  
772 empirically; (d) we assume that direct effects come from a two-component Gaussian mixture,  
773 while they allow for larger number of components; (e) their likelihood function does not explicitly  
774 model the shift between univariate vs multivariate effects (i.e. the LD); (f) CAUSE adds a prior  
775 distribution for the causal/confounder effects and the proportion  $\pi_u$ , while LHC-MR does not;  
776 (g) to calculate the significance of the causal effect they estimate the difference in the expected  
777 log point-wise posterior density and its variance through importance sampling, whereas we use  
778 a simple block jackknife method. Because of point (a), the CAUSE model can be viewed as a  
779 special case of ours when there is no reverse causal effect. We have the advantage of fitting all  
780 parameters simultaneously, while they only approximate this procedure. Although they allow  
781 for more than a two-component Gaussian mixture, for most traits with realistic sample sizes we  
782 do not have enough power to distinguish whether two or more components fit the data better.  
783 Therefore, we believe that a two component Gaussian is a reasonable simplification. Due to the  
784 more complicated approach described in points (e-g), CAUSE is computationally more intense  
785 than LHC-MR, taking up to 1.25 CPU-hours in contrast to our 2.5 CPU-minute run time for a  
786 single starting point optimisation (which is massively parallelisable).

787 When we compared the performance of CAUSE and LHC-MR, we found that for large sample  
788 sizes both LHC-MR and CAUSE performed well not only when applied to data simulated by  
789 their own model, but also by the model of the other method. For smaller sample sizes, both  
790 methods performed poorly when applied to data generated by the other model. However, LHC-  
791 MR was less biased when applied to data generated by its own model than CAUSE was on  
792 data simulated based on its own model, where it provided rather conservative estimates. This  
793 is somewhat expected, since the primary aim of CAUSE is model selection and it is less geared  
794 towards parameter estimation, especially for settings where both sharing and causal effects are  
795 present (leading to very broad estimates). Also, CAUSE parameter estimates have shown to be  
796 somewhat sensitive to the choice of the prior.

797 Finally, when applying both LHC-MR and CAUSE to 156 complex trait pairs, we observed  
798 that the causal effects are reasonably well correlated (0.38 for all estimates, 0.59 for significant  
799 estimates) and agree in sign for trait pairs deemed significantly causal by either or both methods.  
800 In addition, LHC-MR causal estimates were more similar to those of IVW than the estimates  
801 provided by CAUSE. Surprisingly, when a confounding factor was identified by both methods,  
802 the confounder effects (LHC-MR  $t_y/t_x$  ratio and CAUSE  $\eta$  parameter) were uncorrelated. There  
803 are two possible explanations for this: (i) CAUSE may confuse/merge the confounder with the  
804 reverse causal effect, since it does not explicitly model the latter one. (ii) The two models assume  
805 different marginal effect size distributions, hence when multiple heterogeneous confounders exist,  
806 one method may detect one of the confounders, while the other method picks up the other  
807 confounder, depending on which has more similar genetic architecture to the assumed one.

808 Our approach has its own limitations, which we list below. Like any MR method, LHC-MR pro-  
809 vides biased causal effect estimates if the input summary statistics are flawed (e.g. not corrected  
810 for complex population stratification, parental/dynasty effects). As mentioned in the Methods

811 section, our model is strictly-speaking unidentifiable and two distinct sets of parameters fit the  
812 data equally well, if the alternate set of parameters fall within the parameter ranges. As opposed  
813 to classical MR methods that give a single (biased) causal effect estimate, ours can detect and  
814 calculate the competing model. Due to biological considerations, from these competing models,  
815 we chose the one which yielded larger direct heritability than confounder-driven (indirect) heri-  
816 tability. Additional pointers to decide which parameter optimum we choose can be to pick the  
817 one with smaller magnitude of causal effects (large causal effects are unrealistic) or pick the one  
818 that includes causal effects that agree better with those of other MR methods.

819 LHC-MR is not an optimal solution for traits whose genetic architecture substantially deviates  
820 from a two-component Gaussian mixture of effect sizes. Also, for traits with low heritability  
821 ( $< 2.5\%$ ), it is particularly important to compare the causal effect estimates to those from  
822 standard MR methods as results from LHC-MR may be less robust. In addition, trait pairs  
823 with multiple confounders with heterogeneous effect ratios can violate the single confounder  
824 assumption of the LHC model and can lead to biased causal effect estimates. Finally, LHC-MR,  
825 like other methods, is not immune to parental effects that are correlated with offspring effects.  
826 In such cases, the parental effect is grouped with the exposure (due to their strong genetic  
827 correlation) and not viewed as a confounder of the exposure-outcome relationship.

## 828 Acknowledgements

829 This research has been conducted using the UK Biobank Resource under Application Num-  
830 ber 16389. LD scores were calculated based on the UK10K data resource (EGAD00001000740,  
831 EGAD00001000741). Z.K. was funded by the Swiss National Science Foundation (31003A\_169929,  
832 310030\_189147 and 32003B\_173092). For computations we used the CHUV HPC cluster.  
833 We would like to thank Jack Bowden, Valentin Rousson, Matthew Robinson, George Davey  
834 Smith, Thomas Richardson and Eleonora Porcu for their valuable feedback and comments on  
835 this manuscript.

## 836 Author contributions

837 Z.K. devised and directed the project. Z.K., N.M., and L.D. contributed to the mathematical  
838 derivations, design and implementation of the research, to the analysis of the results and to the  
839 writing of the manuscript.

## 840 Data availability

841 The origin of the summary statistics data used is referenced in Table [S1](#). UK Biobank sum-  
842 mary statistics can be downloaded from <http://www.nealelab.is/uk-biobank>. CAD GWAS  
843 summary statistics data from <http://www.cardiogramplusc4d.org/data-downloads/>

## 844 Code availability

845 The source code for this work can be found on [https://github.com/LizaDarrous/LHC-MR\\_](https://github.com/LizaDarrous/LHC-MR_v2/)  
846 [v2/](#)

## 847 References

- 848 [1] Fewell, Z., Davey Smith, G., and Sterne, J. A. C. (2007). The impact of residual and  
849 unmeasured confounding in epidemiologic studies: a simulation study. *American journal of*  
850 *epidemiology* *166*, 646–655.
- 851 [2] Pingault, J.-B., O’Reilly, P. F., Schoeler, T., Ploubidis, G. B., Rijdsdijk, F., and Dudbridge,  
852 F. (2018). Using genetic data to strengthen causal inference in observational research.  
853 *Nature reviews. Genetics* *19*, 566–580.
- 854 [3] Barter, P. J., Caulfield, M., Eriksson, M., Grundy, S. M., Kastelein, J. J. P., Komajda,  
855 M., Lopez-Sendon, J., Mosca, L., Tardif, J.-C., Waters, D. D., et al. (2007). Effects  
856 of torcetrapib in patients at high risk for coronary events. *The New England journal of*  
857 *medicine* *357*, 2109–2122.
- 858 [4] Fordyce, C. B., Roe, M. T., Ahmad, T., Libby, P., Borer, J. S., Hiatt, W. R., Bristow,  
859 M. R., Packer, M., Wasserman, S. M., Braunstein, N., et al. (2015). Cardiovascular drug  
860 development: is it dead or just hibernating? *Journal of the American College of Cardiology*  
861 *65*, 1567–1582.
- 862 [5] Lawlor, D. A., Harbord, R. M., Sterne, J. A. C., Timpson, N., and Davey Smith, G.  
863 (2008). Mendelian randomization: using genes as instruments for making causal inferences  
864 in epidemiology. *Statistics in medicine* *27*, 1133–1163.
- 865 [6] Burgess, S., Butterworth, A., and Thompson, S. G. (2013). Mendelian randomization  
866 analysis with multiple genetic variants using summarized data. *Genetic epidemiology* *37*,  
867 658–665.
- 868 [7] Bowden, J., Davey Smith, G., and Burgess, S. (2015). Mendelian randomization with invalid  
869 instruments: effect estimation and bias detection through egger regression. *International*  
870 *journal of epidemiology* *44*, 512–525.
- 871 [8] Bowden, J., Davey Smith, G., Haycock, P. C., and Burgess, S. (2016). Consistent estima-  
872 tion in mendelian randomization with some invalid instruments using a weighted median  
873 estimator. *Genetic epidemiology* *40*, 304–314.
- 874 [9] Hartwig, F. P., Davey Smith, G., and Bowden, J. (2017). Robust inference in summary  
875 data mendelian randomization via the zero modal pleiotropy assumption. *International*  
876 *journal of epidemiology* *46*, 1985–1998.
- 877 [10] Morrison, J., Knoblauch, N., Marcus, J. H., Stephens, M., and He, X. (2020). Mendelian  
878 randomization accounting for correlated and uncorrelated pleiotropic effects using genome-  
879 wide summary statistics. *Nature Genetics* *52*, 740–747.
- 880 [11] Zhao, Q., Wang, J., Hemani, G., Bowden, J., and Small, D. S. (2020). Statistical inference  
881 in two-sample summary-data Mendelian randomization using robust adjusted profile score.  
882 *The Annals of Statistics* *48*, 1742 – 1769.
- 883 [12] Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Schizophrenia Work-  
884 ing Group of the Psychiatric Genomics Consortium, Patterson, N., Daly, M. J., Price, A. L.,  
885 and Neale, B. M. (2015). Ld score regression distinguishes confounding from polygenicity  
886 in genome-wide association studies. *Nature genetics* *47*, 291–295.
- 887 [13] Jordan, D. M., Verbanck, M., and Do, R. (2019). The landscape of pervasive horizontal  
888 pleiotropy in human genetic variation is driven by extreme polygenicity of human traits  
889 and diseases. *bioRxiv*.

- 890 [14] Bulik-Sullivan, B., Finucane, H. K., Anttila, V., Gusev, A., Day, F. R., Loh, P.-R., Re-  
891 proGen Consortium, Psychiatric Genomics Consortium, Genetic Consortium for Anorexia  
892 Nervosa of the Wellcome Trust Case Control Consortium 3, Duncan, L., et al. (2015). An  
893 atlas of genetic correlations across human diseases and traits. *Nature genetics* *47*, 1236–  
894 1241.
- 895 [15] Nadarajah, S. and Pogány, T. K. (2016). On the distribution of the product of correlated  
896 normal random variables. *Comptes Rendus Mathématique* *354*, 201–204.
- 897 [16] McNolty, F. (1973). Some probability density functions and their characteristic functions.
- 898 [17] Bateman, H. (1953). Volume i.
- 899 [18] Heideman, M. T., Johnson, D. H., and Burrus, C. S. (1985). Gauss and the history of the  
900 fast fourier transform. *Archive for History of Exact Sciences* *34*, 265–277.
- 901 [19] Neale Lab (2018). UK BioBank. <http://www.nealelab.is/uk-biobank/>.
- 902 [20] Walter, K., Min, J. L., Huang, J., Crooks, L., Memari, Y., McCarthy, S., Perry, J. R. B.,  
903 Xu, C., Futema, M., Lawson, D., et al. (2015). The uk10k project identifies rare variants  
904 in health and disease. *Nature* *526*, 82–90.
- 905 [21] Rüeger, S., McDaid, A., and Kutalik, Z. (2018). Evaluation and application of summary  
906 statistic imputation to discover new height-associated loci. *PLoS genetics* *14*, e1007371.
- 907 [22] 1000 Genomes Project Consortium. (2010). A map of human genome variation from  
908 population-scale sequencing. *Nature* *467*, 1061–1073.
- 909 [23] R Core Team. (2019). R: A Language and Environment for Statistical Computing. R  
910 Foundation for Statistical Computing Vienna, Austria.
- 911 [24] MRC IEU (2019). EpiGraphDB. <http://epigraphdb.org/>.
- 912 [25] Liu, Y., Elsworth, B., Erola, P., Haberland, V., Hemani, G., Lyon, M., Zheng, J., and  
913 Gaunt, T. R. (2020). Epigraphdb: A database and data mining platform for health data  
914 science. *bioRxiv*.
- 915 [26] Clogg, C. C., Petkova, E., and Haritou, A. (1995). Statistical methods for comparing  
916 regression coefficients between models. *American Journal of Sociology* *100*, 1261–1293.
- 917 [27] Hemani, G., Zheng, J., Elsworth, B., Wade, K. H., Haberland, V., Baird, D., Laurin, C.,  
918 Burgess, S., Bowden, J., Langdon, R., et al. (2018). The MR-Base platform supports  
919 systematic causal inference across the human phenome. *eLife* *7*, e34408.
- 920 [28] Allen, M., Poggiali, D., Whitaker, K., Marshall, T., and Kievit, R. (2019). Raincloud plots:  
921 a multi-platform tool for robust data visualization [version 1; peer review: 2 approved].  
922 *Wellcome Open Research* *4*.
- 923 [29] Laclaustra, M., Lopez-Garcia, E., Civeira, F., Garcia-Esquinas, E., Graciani, A., Guallar-  
924 Castillon, P., Banegas, J. R., and Rodriguez-Artalejo, F. (2018). Ldl cholesterol rises with  
925 bmi only in lean individuals: Cross-sectional u.s. and spanish representative data. *Diabetes*  
926 *Care* *41*, 2195–2201.
- 927 [30] Drøyvold, W. B., Midthjell, K., Nilsen, T. I. L., and Holmen, J. (2005). Change in body  
928 mass index and its impact on blood pressure: a prospective population study. *International*  
929 *Journal of Obesity* *29*, 650–655.

- 930 [31] Lee, M.-R., Lim, Y.-H., and Hong, Y.-C. (2018). Causal association of body mass index with  
931 hypertension using a mendelian randomization design. *Medicine (Baltimore)* *97*, e11252.
- 932 [32] Corbin, L. J., Richmond, R. C., Wade, K. H., Burgess, S., Bowden, J., Smith, G. D., and  
933 Timpson, N. J. (2016). Bmi as a modifiable risk factor for type 2 diabetes: Refining and  
934 understanding causal estimates using mendelian randomization. *Diabetes* *65*, 3002–3007.
- 935 [33] Narayan, K., Boyle, J. P., Thompson, T. J., Gregg, E. W., and Williamson, D. F. (2007).  
936 Effect of bmi on lifetime risk for diabetes in the u.s. *Diabetes Care* *30*, 1562–1566.
- 937 [34] Yusuf, S., Hawken, S., Ounpuu, S., Bautista, L., Franzosi, M. G., Commerford, P., Lang,  
938 C. C., Rumboldt, Z., Onen, C. L., Lisheng, L., et al. (2005). Obesity and the risk of  
939 myocardial infarction in 27,000 participants from 52 countries: a case-control study. *Lancet*  
940 *366*, 1640–1649.
- 941 [35] Riaz, H., Khan, M. S., Siddiqi, T. J., Usman, M. S., Shah, N., Goyal, A., Khan, S. S.,  
942 Mookadam, F., Krasuski, R. A., and Ahmed, H. (2018). Association Between Obesity and  
943 Cardiovascular Outcomes: A Systematic Review and Meta-analysis of Mendelian Random-  
944 ization Studies. *JAMA Network Open* *1*, e183788–e183788.
- 945 [36] Sun, D., Zhou, T., Heianza, Y., Li, X., Fan, M., Fonseca, V. A., and Qi, L. (2019). Type  
946 2 diabetes and hypertension. *Circulation research* *124*, 930–937.
- 947 [37] Tomeo, C. A., Field, A. E., Berkey, C. S., Colditz, G. A., and Frazier, A. L. (1999). Weight  
948 concerns, weight control behaviors, and smoking initiation.
- 949 [38] Cawley, J., Markowitz, S., and Tauras, J. (2004). Lighting up and slimming down: the  
950 effects of body weight and cigarette prices on adolescent smoking initiation.
- 951 [39] Cao, M. and Cui, B. (2020). Association of educational attainment with adiposity, type  
952 2 diabetes, and coronary artery diseases: A mendelian randomization study. *Frontiers in*  
953 *Public Health* *8*, 112.
- 954 [40] Loucks, E. B., Buka, S. L., Rogers, M. L., Liu, T., Kawachi, I., Kubzansky, L. D., Martin,  
955 L. T., and Gilman, S. E. (2012). Education and coronary heart disease risk associations  
956 may be affected by early-life common prior causes: a propensity matching analysis. *Ann*  
957 *Epidemiol* *22*, 221–232.
- 958 [41] Gage, S. H., Bowden, J., Davey Smith, G., and Munafò, M. R. (2018). Investigating causal-  
959 ity in associations between education and smoking: a two-sample Mendelian randomization  
960 study. *International Journal of Epidemiology* *47*, 1131–1140.
- 961 [42] Sanderson, E., Davey Smith, G., Bowden, J., and Munafò, M. R. (2019). Mendelian  
962 randomisation analysis of the effect of educational attainment and cognitive ability on  
963 smoking behaviour. *Nature Communications* *10*, 2949.
- 964 [43] Marouli, E., Del Greco, M. F., Astley, C. M., Yang, J., Ahmad, S., Berndt, S. I., Caulfield,  
965 M. J., Evangelou, E., McKnight, B., Medina-Gomez, C., et al. (2019). Mendelian ran-  
966 domisation analyses find pulmonary factors mediate the effect of height on coronary artery  
967 disease. *Communications biology* *2*, 119.
- 968 [44] Tan, L. E., Llano, A., Aman, A., Dominiczak, A. F., and Padmanabhan, S. (2018). A18709  
969 mendelian randomization study of causal relationship of height on blood pressure and ar-  
970 terial stiffness. *Journal of Hypertension* *36*.

- 971 [45] Laaksonen, D. E., Niskanen, L., Nyysönen, K., Lakka, T. A., Laukkanen, J. A., and  
972 Salonen, J. T. (2008). Dyslipidaemia as a predictor of hypertension in middle-aged men.  
973 *Eur Heart J* *29*, 2561–2568.
- 974 [46] Davies, N. M., Howe, L. J., Brumpton, B., Havdahl, A., Evans, D. M., and Davey Smith,  
975 G. (2019). Within family Mendelian randomization studies. *Human Molecular Genetics*  
976 *28*, R170–R179.
- 977 [47] Benson, R., von Hippel, P. T., and Lynch, J. L. (2018). Does more education cause  
978 lower bmi, or do lower-bmi individuals become more educated? evidence from the national  
979 longitudinal survey of youth 1979. *Soc Sci Med* *211*, 370–377.
- 980 [48] Witter, F. R. and Luke, B. (1991). The effect of maternal height on birth weight and birth  
981 length. *Early Human Development* *25*, 181–186.
- 982 [49] Tyrrell, J., Jones, S. E., Beaumont, R., Astley, C. M., Lovell, R., Yaghootkar, H., Tuke,  
983 M., Ruth, K. S., Freathy, R. M., Hirschhorn, J. N., et al. (2016). Height, body mass index,  
984 and socioeconomic status: mendelian randomisation study in uk biobank. *BMJ* *352*.
- 985 [50] Bulik-Sullivan, B., Finucane, H. K., Anttila, V., Gusev, A., Day, F. R., Loh, P.-R., Duncan,  
986 L., Perry, J. R. B., Patterson, N., Robinson, E. B., et al. (2015). An atlas of genetic  
987 correlations across human diseases and traits. *Nature Genetics* *47*, 1236–1241.
- 988 [51] Mounier, N. and Kutalik, Z. (2021). Correction for sample overlap, winner’s curse and  
989 weak instrument bias in two-sample mendelian randomization. *bioRxiv*.
- 990 [52] Thomas, D., Strauss, J., and Henriques, M.-H. (1991). How does mother’s education affect  
991 child height? *The Journal of Human Resources* *26*, 183–211.
- 992 [53] Warrington, N. M., Beaumont, R. N., Horikoshi, M., Day, F. R., Helgeland, Ø., Laurin,  
993 C., Bacelis, J., Peng, S., Hao, K., Feenstra, B., et al. (2019). Maternal and fetal genetic  
994 effects on birth weight and their relevance to cardio-metabolic risk factors. *Nat Genet* *51*,  
995 804–814.
- 996 [54] Brumpton, B., Sanderson, E., Hartwig, F. P., Harrison, S., Vie, G. Å., Cho, Y., Howe, L. D.,  
997 Hughes, A., Boomsma, D. I., Havdahl, A., et al. (2019). Within-family studies for mendelian  
998 randomization: avoiding dynastic, assortative mating, and population stratification biases.  
999 *bioRxiv*.
- 1000 [55] O’Connor, L. J. and Price, A. L. (2018). Distinguishing genetic correlation from causation  
1001 across 52 diseases and complex traits. *Nature genetics* *50*, 1728–1734.
- 1002 [56] Brown, B. C. and Knowles, D. A. (2020). Phenome-scale causal network discovery with  
1003 bidirectional mediated mendelian randomization. *bioRxiv*.
- 1004 [57] Howey, R., Shin, S.-Y., Relton, C., Smith, G. D., and Cordell, H. J. (2019). Bayesian  
1005 network analysis incorporating genetic anchors complements conventional mendelian ran-  
1006 domization approaches for exploratory analysis of causal relationships in complex data.  
1007 *bioRxiv*.