# Mendelian randomization for multiple exposures and outcomes with Bayesian Directed Acyclic Graphs exploration and causal effects estimation

Verena Zuber[1,2,3,*], Toinét Cronjé[4], Na Cai[5,6,7], Dipender Gill[1], and Leonardo Bottolo[8,9,*]

[1]Department of Epidemiology and Biostatistics, School of Public Health, Imperial College London, London, UK
[2]MRC Centre for Environment and Health, School of Public Health, Imperial College London, London, UK
[3]UK Dementia Research Institute, Imperial College London, London, UK
[4]Department of Public Health, University of Copenhagen, Denmark
[5]Helmholtz Pioneer Campus, Helmholtz Munich, Neuherberg, Germany
[6]Computational Health Centre, Helmholtz Munich, Neuherberg, Germany
[7]School of Medicine and Health, Technical University of Munich, Munich, Germany
[8]Department of Medical Genetics, School of Clinical Medicine, University of Cambridge, Cambridge, UK
[9]MRC Biostatistics Unit, School of Clinical Medicine, University of Cambridge, Cambridge, UK
[*]Corresponding authors: Verena Zuber v.zuber@imperial.ac.uk and Leonardo Bottolo lb664@cam.ac.uk

**Abstract**

Current Mendelian randomization (MR) methods do not reflect complex relationships among multiple exposures and outcomes as is typical for real-life applications. We introduce MrDAG the first MR method to model dependency relations within the exposures, the outcomes, and between them to improve causal effects estimation. MrDAG combines three causal inference strategies in a unified manner. It uses genetic variation as instrumental variables to account for unmeasured confounders. It performs structure learning to detect and orientate the direction of the dependencies within exposures and outcomes. Finally, interventional calculus is employed to derive principled causal effect estimates. MrDAG was motivated to unravel how lifestyle and behavioural exposures impact mental health. It highlights education and smoking as key effective points of intervention given their downstream effects on mental health. These insights would have been difficult to delineate without modelling the causal paths between multiple exposures and outcomes at once.

# Introduction

Genetic evidence is increasingly used to infer causal relationships between human traits in Mendelian randomization (MR) analysis. The standard MR paradigm, one exposure and one outcome, can be biased by unobserved pleiotropy. It occurs when the genetic variants used as instruments in the MR analysis act *via* separate pathways to the exposure under investigation. Extensions to consider multiple exposures [1] along with multi-response [2] of standard MR allow to model pleiotropy acting either *via* any of the exposures or any of the outcomes or both, respectively.

Yet, to date, there is no MR approach which can estimate the dependency relations within the exposures and the outcomes to enhance the detection of causal effects between them and improve their accuracy. As we show in our motivating data application on mental health phenotypes, it is a common problem in practical applications that the effect of an exposure on an outcome can be confounded or (partially or completely) mediated by another exposure [3] or mediated by another outcome, or both. However, this structure is latent and not known and consequently needs to be learned from the data. This problem has been overlooked in the literature and current MR implementations which do not account for these dependencies likely produce spurious findings which are often claimed as supporting causality in applied analysis.

Here, we address this gap by proposing the MrDAG model, the first Mendelian randomization method with Directed Acyclic Graphs (DAGs) exploration and causal effects estimation, which utilises summary-level genetic associations from genome-wide association studies to learn how interrelated exposures affect multiple outcomes which, in turn, are interconnected in a complex fashion. MrDAG is a Bayesian causal graphical model that combines three causal inference strategies in a unified manner. First, the MR paradigm which uses genetic variation as instrumental variables (IVs) [4, 5] to ensure unconfoundedness. Second, structure learning [6], *i.e.*, graphical models selection to define the graphs that best describe the dependency structure in a given data set under the constraint on edges' orientation from the exposures to the outcomes implied by the MR paradigm. Third, interventional calculus to derive principled causal effects estimates [7] from the exposures to the outcomes.

Our motivating real data application considers the impact of six common modifiable lifestyle and behavioural exposures on seven mental health phenotypes. Mental health describes patterns of cognitive, emotional, and behavioural disregulations that limit daily functioning and cause distress. One in eight individuals suffers from one or more mental health phenotypes worldwide, most commonly anxiety, attention-deficit hyperactivity, autism spectrum, bipolar, eating, personality or schizophrenia-related diseases [8]. Collectively, they contribute to more than 15% of total years lived with disability [9]. Clinically, mental health phenotypes are notoriously difficult to disentangle and diagnose due to

the lack of objective biological biomarkers and distinct disease impressions [10]. No symptom can be uniquely ascribed to one disease, and each disease comprises experiencing a group of interrelated traits. In research, this complexity is reinforced by the multifaceted mechanisms that cause and sustain mental health [10, 11]. In addition to genetic liability, numerous behavioural and lifestyle factors such as alcohol consumption, smoking, sleep hygiene, physical activity and education contribute to the risk of developing a mental health trait [11, 12]. Notably, these factors are also affected by existing disease and treatment [13]. It is essential to appreciate these complexities when attempting to identify distinct and shared underlying mechanisms of mental health. While MR studies have been effective in circumventing some of the limitations of traditional epidemiology such as environmental confounding and reverse causation, MR remains largely unable to fully disentangle the interplay between traits that cause or result from mental health [14]. The complexity of such an example demonstrates the limitations of current MR solutions to offer a more comprehensive picture of causal mechanisms between complex phenotypes and provides a suitable test ground for the application of the proposed methodology.
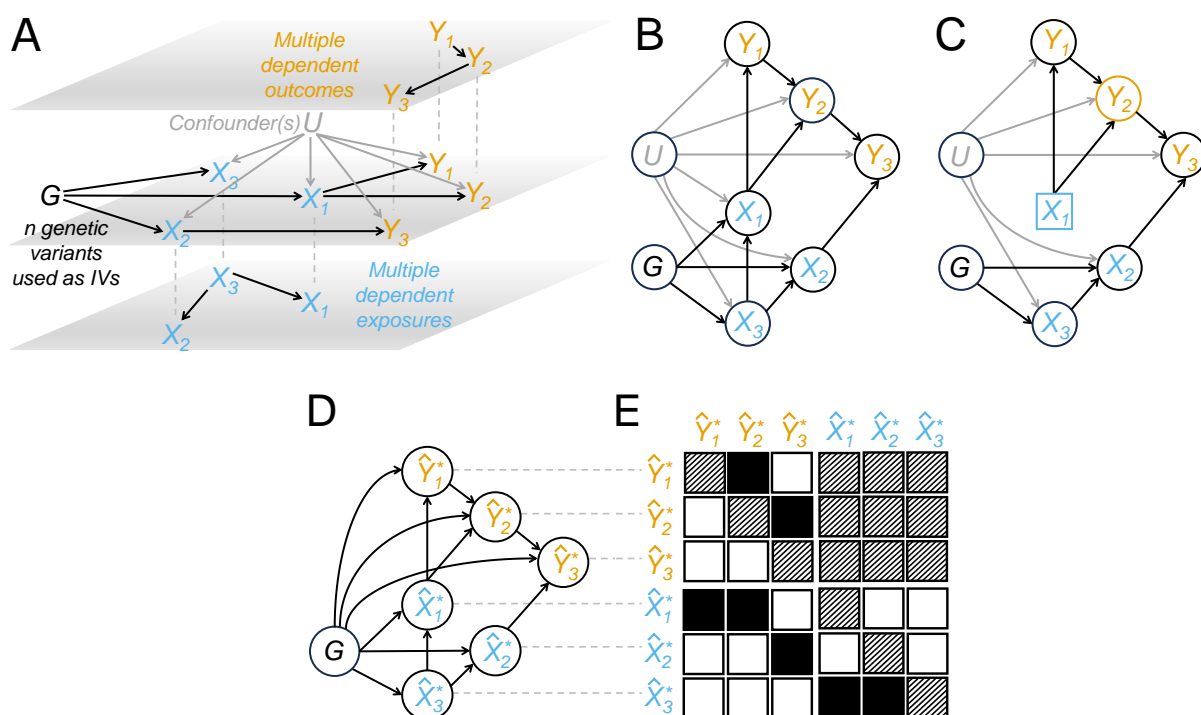
# Results

## Causal inferential strategies in MrDAG

MrDAG combines three causal inference strategies.

First, MR has pioneered the ability to use genetic data as IVs to derive causal statements from observational data despite the presence of unmeasured confounders [15, 16].

Second, in its standard formulation of one exposure and one outcome, the conditional dependencies between the outcome $Y$, the exposure $X$, the IV $G$ and the unmeasured confounder $U$ are all given as well as their graphical representation [5]. When multiple exposures $\boldsymbol{X}$ [1] and multiple outcomes $\boldsymbol{Y}$ [2] are considered along with multiple IVs $\boldsymbol{G}$, (partial) correlation between $\boldsymbol{X}$ and conditional dependencies between $\boldsymbol{Y}$ are included in the models to perform the selection of important exposures whose causal effects can be shared or are distinct across the responses. However, to date, no dependency relations within the exposures and within the outcomes are estimated, although, in practical applications, the effect of an exposure on an outcome can be confounded or (partially or completely) mediated by another exposure or mediated by another outcome, or both, see Figures 1A-B for an illustration. Moreover, dependency relations are also important to derive principled estimates of the causal effects [7].

In real data applications, complex dependency relations between the traits are generally not known in advance, and they need to be learned from the data. To detect them, we rely on Directed Acyclic Graphs (DAGs) and structure learning. Graphical models

**Figure 1. Directed Acyclic Graph (DAG) representation of the proposed multiple exposures and multiple outcomes Mendelian randomization model and causal effects estimation.** (**A**) Middle panel: Multivariable Mendelian randomization for multiple responses with $\boldsymbol{G} = (G_1, \ldots, G_n)^\top$: Genetic variants (black) or instrumental variables (IVs); $\boldsymbol{X} = (X_1, X_2, X_3)^\top$: Exposures (blue); $\boldsymbol{Y} = (Y_1, Y_2, Y_3)^\top$: Responses (orange); $U$: Unmeasured confounder(s) (grey). True (unconfounded by $U$) exposure-outcome dependency relations depicted in the middle panel are classified as follows: $X_1$ has *shared* causal effect on responses $Y_1$ and $Y_2$, while $X_2$ has a *distinct* causal effect on response $Y_3$. $X_3$ does not have any effect on the outcomes. Bottom panel: True fork structure within the exposures with $X_3$ regarded as the common cause of $X_1$ and $X_2$. Top panel: True chain structure within the outcomes, where $Y_1$ affects $Y_3$ through $Y_2$. (**B**) DAG is obtained by combining the true exposure-outcome dependency relations ((A) middle panel), the fork structure within the exposures ((A) bottom panel) and the chain structure within the responses ((A) top panel). When looking at the effect of $X_1$ on $Y_3$, $X_3$ (along with $\boldsymbol{G}$ and $U$) is a confounder of $X_1$ and $Y_2$ is a complete mediator. Without conditioning on $Y_2$, with the same set of confounders, a spurious association would be found between $X_1$ and $Y_3$. (**C**) Estimation of the causal effect under intervention in $X_1$ on $Y_2$, highlighted in blue and orange, respectively. The representation of $X_1$ has changed from a circle to a square to emphasise that, under intervention, it is no longer a random variable and it is now set at $X_1 = \widetilde{x}_1$. Intervention affects only the conditional distribution of $X_1$, *i.e.*, $X_1 \mid (X_3, \boldsymbol{G}, U)$ and it leaves unaltered all the others. From a practical perspective, it would be sufficient to condition on $X_3$, $\boldsymbol{G}$ and $U$ (graphically, the directed edges to $X_1$ from $X_3$, $\boldsymbol{G}$ and $U$ are removed) to guarantee that the association between $X_1$ and $Y_2$ is purely causative (see Supplementary Figure 1). However, since $U$ is unobserved, the estimation of the causal effects cannot be obtained only by conditioning. (**D**) Genetically predicted exposures $\widehat{\boldsymbol{X}}^* = (\widehat{X}_1^*, \widehat{X}_2^*, \widehat{X}_3^*)^\top$ and outcomes $\widehat{\boldsymbol{Y}}^* = (\widehat{Y}_1^*, \widehat{Y}_2^*, \widehat{Y}_3^*)^\top$ depend only on $\boldsymbol{G}$ which are chosen to be associated only with $\boldsymbol{X}$ and not with $\boldsymbol{Y}$. Graphically, no directed edges to $\widehat{\boldsymbol{X}}^*$ and $\widehat{\boldsymbol{Y}}^*$ from $U$ are pictured. True (unconfounded by $U$) dependency relations between the traits in the original (individual-level) data shown in (B) are obtained by the DAG estimated by using $\widehat{\boldsymbol{X}}^*$ and $\widehat{\boldsymbol{Y}}^*$. (**E**) Adjacency matrix describing the Markov properties of the DAG obtained by using genetically predicted exposures and outcomes (the variables in the $x$-axis are dependent on the variables in the $y$-axis) which are function of the IVs and the inverse-variance weighted (IVW) (depicted with a " * ") summary-level statistics $\widehat{\boldsymbol{B}}_X^* = (\widehat{\boldsymbol{\beta}}_{X_1}^*, \widehat{\boldsymbol{\beta}}_{X_2}^*, \widehat{\boldsymbol{\beta}}_{X_3}^*)^\top$ and $\widehat{\boldsymbol{B}}_Y^* = (\widehat{\boldsymbol{\beta}}_{Y_1}^*, \widehat{\boldsymbol{\beta}}_{Y_2}^*, \widehat{\boldsymbol{\beta}}_{Y_3}^*)^\top$. Neither reverse causation (top-right submatrix) nor feedback loops (main diagonal) are allowed. Colour code: Black, directed edge between variables; white, no causal relationship between variables; black-white strips, directed edge not allowed (feedback loop and reverse causation between exposures and outcomes).

are multivariate distributions associated with a graph and are very effective for encoding conditional dependencies [17] between random variables. They are represented in a graph as nodes (vertices) while edges denote conditional dependence relationships between the corresponding random variables. A DAG is a directed graph, where each edge has an orientation with no directed cycles. Structure learning is a model selection problem [6] to estimate the DAG (or competing DAGs) that best describes the dependency structure in a given data set. However, without identifiability conditions [18], it is not possible to estimate uniquely the underlying DAG since its conditional independencies can be associated with several alternative DAGs. The set of DAGs that hold the same conditional independencies is known as Markov Equivalent Class and the best that can be done from observational data is to estimate this class (or competing classes). Thus, this paper aims to illustrate how to perform DAG exploration (whose importance will be apparent in the next paragraph) which belongs to the Markov Equivalent Classes that best fit the data under the constraint on the orientation of the edges, known as partial ordering [19], from the exposures to the outcomes implied by the MR paradigm.

Third, besides the identification of the exposure-outcome relations as well as the dependency patterns within the exposures and the outcomes, we are also interested in the causal effects estimation under intervention [7]. An intervention on the exposures can be made explicit by a suitable modification of the multivariate distribution associated with the DAG, under the assumption that the intervention does not affect any other variable in the joint distribution besides the conditional distribution of the exposure under intervention [20]. See Figure 1C for an example of intervention on an exposure and the estimation of the causal effect on an outcome.

In this formulation, all confounders should be measurable to perform structure learning and causal effects estimation (causal sufficiency assumption [21]). This assumption is usually not met in real data applications where, instead, unmeasured confounders are ubiquitous and affect exposures and responses at the same time. To solve this problem, we demonstrate (see Methods) and show in an extensive simulation study (see Results) that, under partial ordering, we can estimate the dependency structure that exists between the traits in the original (individual-level) data unconfounded by $U$ by using their genetically predicted values. Since the genetically predicted traits depend only on the selected IVs, the confounders do not mask the true dependency relations required in causal effects estimation. See Figure 1D, where the graphical model estimated by using genetically predicted exposures and outcomes approximates the corresponding DAG in the individual-level data not affected by $U$. Our approach shares some similarities with methods based on the genetic correlation and developed to analyse the joint genetic architecture of complex traits [22] although, in the proposed MR framework, genetic variants are chosen to be valid IVs in contrast to genetic variants chosen for genome-wide [23] or local genetic correlation [24]. Computationally, given the duality between the Markov properties of the DAG and

a non-symmetrical adjacency matrix (see Figure 1E), structure learning of the graphical model (or competing graphical models) that best fits the data is performed on a non-symmetrical adjacency matrix which incorporates the constraints on the orientation of the edges from the exposures to the outcomes.

Finally, for a given DAG, we extend results regarding the consistency of the effects of the regressions of the exposures and the outcomes on $\boldsymbol{G}$ which can be obtained without adjustment on $U$ since the genetic variants used as IVs are randomly assigned [4], and show that it is possible to identify and estimate the causal effects between multiple exposures and multiple outcomes based on Pearl's interventional calculus [7] (see Methods).

The MrDAG model can be summarised as follows:

$$[\boldsymbol{g}^\top \widehat{\boldsymbol{B}}_Y^* \boldsymbol{g}^\top \widehat{\boldsymbol{B}}_X^*]^\top \sim \mathrm{N}_{q+p}([\boldsymbol{g}^\top \boldsymbol{B}_Y^* \boldsymbol{g}^\top \boldsymbol{B}_X^*]^\top, \boldsymbol{\Sigma}^*),$$

where $\boldsymbol{g}$ are the observed IVs after pruning or clumping, $\widehat{\boldsymbol{B}}_Y^*$ and $\widehat{\boldsymbol{B}}_X^*$ are the inverse-variance weighted (IVW) [25] estimated genetic associations with the outcomes and the exposures, $\boldsymbol{g}^\top \widehat{\boldsymbol{B}}_Y^*$ and $\boldsymbol{g}^\top \widehat{\boldsymbol{B}}_X^*$ are the genetically predicted values of the outcomes $\widehat{\boldsymbol{Y}}^*$ and exposures $\widehat{\boldsymbol{X}}^*$ based on the IVs, respectively, which are normally distributed for large sample sizes, and $\boldsymbol{\Sigma}^*$ is the genetic covariance matrix that can be partitioned into $\boldsymbol{\Sigma}_{XX}^*$, $\boldsymbol{\Sigma}_{YY}^*$ and $\boldsymbol{\Sigma}_{XY}^*$, the genetic covariances within the exposures, the outcomes and between them. The MrDAG model allows us to find a solution to the two problems highlighted before. First, we perform DAGs exploration under partial ordering by using $\boldsymbol{\Omega}^* = \boldsymbol{\Sigma}^{*-1}$, to learn the unconfounded dependency relations within the exposures, the outcomes and between them and to understand the genetic paths that link exposures and outcomes (see Methods). Second, estimate the causal effects of the intervention on the exposures as a function of trait-specific elements of the genetic associations $\widehat{\boldsymbol{B}}_Y^*$ and $\widehat{\boldsymbol{B}}_X^*$ informed by the explored DAGs, unconfounded by any pleiotropic effects within the exposures and the outcomes and any unmeasured confounder.

## Selection of instrumental variables

MrDAG uses the same instrument selection employed in MVMR regardless of the multiple outcomes [2]. A genetic variant is considered a valid instrument for MVMR when three core conditions hold [3]: (IV1) Independence: The variant is independent of all confounders of each of the exposure-outcome associations; (IV2) Relevance: The variant must not be conditional independent of each exposure given the other exposures; (IV3) Exclusion restriction: The variant is independent of the outcome conditional on the exposures and confounders. In practice, only IV2 can be computationally evaluated from the available data. A recent solution to mitigate the effects of weak IVs in MVMR is presented in [26].

There is an important distinction between IV selection in MVMR, as used by MrDAG,

and bidirectional MR. Let's consider two traits A and B. In bidirectional MR, two MR analyses are conducted, one for trait $A$ on trait $B$, and then *vice versa*. First, specific IVs are selected for trait A and the first MR model is fit. Then, another set of specific IVs is selected for trait B and the second MR model tests the opposite effects direction. In contrast, in MVMR, IVs are chosen to be the union of genome-wide significant genetic variants for any exposure. By combining MVMR IVs selection approach with DAG learning, MrDAG can infer the bidirectionality of the relationships within exposures based on $\boldsymbol{\Omega}_{XX}^* = \boldsymbol{\Sigma}_{XX}^{*-1}$ without repeated IVs selection and subsequent analyses. A similar comment can be made for the estimation of the bidirectionality of the relationships within the outcomes based on $\boldsymbol{\Omega}_{YY}^* = (\boldsymbol{\Sigma}_{YY}^* - \boldsymbol{\Sigma}_{YX}^* \boldsymbol{\Sigma}_{XX}^{*-1} \boldsymbol{\Sigma}_{XY}^*)^{-1}$ (see Methods). These dependencies should be interpreted as an indication of a violation of condition IV3, *i.e.*, pleiotropy not explained by the estimated causal effects from the exposures to the outcomes [2]. The detected relationships within the exposures also suggest the existence of measured pleiotropy which, in the proposed framework, comprises confounding, mediation and independent pleiotropic pathways [3].

Overall, only the direction from exposures to outcomes is fixed in MrDAG, and no reverse causation is allowed, reflecting the standard MR paradigm.

## Simulation study

We compare MrDAG in an extensive and comprehensive simulation study where four different *in silico* scenarios have been generated on individual-level data for $N = 100,000$ individuals. The simulated data sets include $n = 100$ independent genetic variants $\boldsymbol{G}$, an unmeasured confounder $U$, 15 exposures $\boldsymbol{X}$ and 5 outcomes $\boldsymbol{Y}$. All exposures $\boldsymbol{X}$ were measured on the same individuals in the first sample and have complete overlap as well as all outcomes $\boldsymbol{Y}$ were measured on the same individuals in the second sample independent of the first sample. In all simulations, the unconfounded dependency relations between the traits are simulated at the individual-level while the algorithms use as input the corresponding IVW summary-level statistics.

The four simulation scenarios are built by combining two different strategies we used to simulate the dependency patterns within the exposures and the responses:

i) "UndG$_X$-Med$_Y$". A sparse undirected graphical model ("UndG$_X$") encodes the dependency pattern within the exposures $\boldsymbol{X} = (X_1, \ldots, X_{15})$. Regarding the responses $\boldsymbol{Y} = (Y_1, \ldots, Y_5)$, one outcome is completed mediated by another one ("Med$_Y$"). For a visual representation of this scenario, see Figure 2A.

ii) "DAG$_X$-Med$_Y$". The dependency relations within the exposures are more complex than in scenario (i) since a topologically ordered DAG within the exposures

7

**Figure 2. Schematic illustration of different dependency structures simulated between the traits at the individual-level data and the parameters employed in the simulation study.** Directed edges indicate dependency relations, while undirected edges denote partial correlations. Dashed lines depict the true (unconfounded by $U$) dependency structure within the exposures and the outcomes, while solid lines indicate true causal effects between them. Parameters $\psi_Y$ and $\psi_X$ indicate the simulated effects of the unmeasured confounder $U$ on the exposures and the outcomes, respectively, and $\boldsymbol{B}_X = (\boldsymbol{\beta}_{X_1}, \boldsymbol{\beta}_{X_2}, \boldsymbol{\beta}_{X_3})$ are the simulated genetic effects on the exposures. For simplicity, they are shown only on the left panel. $\boldsymbol{\Theta} = (\theta_{1,1}, \theta_{1,2}, \theta_{2,3})$ are the simulated causal effects from the exposures to the outcomes while $\boldsymbol{\Gamma}_X = (\gamma_{3,1}^X, \gamma_{3,2}^X)$ and $\boldsymbol{\Gamma}_Y = (\gamma_{1,2}^Y, \gamma_{2,3}^Y)$ are the mediation parameters within the exposures and the outcomes, respectively, where the subscripts denote their directionality. When partial correlations are simulated within the exposures, bidirectional effects are depicted with double subscripts, *i.e.*, $\boldsymbol{\Gamma}_X = (\gamma_{1,2//2,1}^X, \gamma_{2,3//3,2}^X)$. (**A**) Simulated scenario "UndG$_X$-Med$_Y$", where an undirected graph ("UndG$_X$") encodes the dependency pattern within $\boldsymbol{X}$ and, within the responses, an outcome ($Y_3$) is completed mediated ("Med$_Y$") by another response ($Y_2$) which, in turn, is affected by a different exposure ($X_1$). Although there is another partial mediation between $X_1$ and $Y_3$ through $X_2$, this mediation happens within $\boldsymbol{X}$, so it does not affect the definition of complete mediation within $\boldsymbol{Y}$. (**B**) Simulated scenario "DAG$_X$-Med$_Y$", where a topologically ordered DAG within the exposures ("DAG$_X$") is simulated. Specifically, in the example depicted, a fork structure is simulated, *i.e.*, $X_3$ affects both $X_1$ and $X_2$. A complete mediation is still considered within the responses. (**C**) Simulated scenario "UndG$_X$-DAG$_Y$". Here, the dependency structure between the individual-level responses is obtained by simulating a topologically ordered DAG ("DAG$_Y$"). Specifically, a chain structure is considered, *i.e.*, $Y_1$ affects $Y_2$ which, in turn, affects $Y_3$, whereas an undirected graph encodes the dependency pattern within $\boldsymbol{X}$. (**D**) Simulated scenario "DAG$_X$-DAG$_Y$", where two topologically ordered DAGs are simulated within the exposures (fork structure) and outcomes (chain structure), respectively.

("DAG$_X$") is simulated [27]. A complete mediation is still considered within the responses. This second scenario is illustrated in Figure 2B.

iii) "UndG$_X$-DAG$_Y$". Here, a more complex dependency structure within the individual-level responses ("DAG$_Y$") is simulated. This scenario is represented in Figure 2C. An example of the complex dependency patterns generated in the simulation study between the traits for one replicate of scenario UndG$_X$-DAG$_Y$ is shown in Figure 3A.

iv) "DAG$_X$-DAG$_Y$": This is the most complex simulated scenario where two independent topologically ordered DAGs have been simulated within the exposures and outcomes. Figure 2D presents a schematic illustration of this scenario, while Figure 3E shows the intricate dependency structure simulated between the traits for one replicate of DAG$_X$-DAG$_Y$ scenario.

Taken together, in scenarios (ii) and (iv), the overall individual-level DAGs, obtained

by combining two different simulation strategies for $X$ and $Y$, are fully oriented while in scenarios (i) and (iii) the overall DAGs are partially oriented. Details regarding the parameters $\psi_X$ and $\psi_Y$, the simulated levels of the effects of the unmeasured confounder $U$ on the responses and the outcomes, $B_X$, the simulated levels of the genetic effects on the exposures, and $\Gamma_X$ and $\Gamma_Y$, the simulated levels of the mediation parameters within the exposures and the outcomes are presented in Methods. Finally, all simulations are replicated 25 times and initialised with a different random seed.

We compare MrDAG with published multivariable MR methods and their software implementations excluding from the comparisons naïve one-exposure and one-outcome MR models since it has been shown that they are outperformed by multivariable MR methods when there is measured pleiotropy among exposures [3]. Specifically, we consider Mendelian randomization with Bayesian Model Averaging (MR-BMA) [1] an MVMR algorithm which allows for many exposures to be included, but does not model explicitly the dependency relations within the exposures [3]. MR-BMA estimates the sparse direct causal effects between the exposures and one outcome providing the marginal posterior probability of inclusion (mPPI) along with the posterior mean of the causal effects. We treat MR-BMA as the baseline algorithm for the comparisons since it analyses one outcome at-a-time. Secondly, we include Mendelian randomization with PC algorithm (MRPC) [28], which combines instrumental variables with the PC algorithm [29] for DAG estimation. At a specified type I error rate for the conditional independence test, MRPC returns the estimated Partially Directed Acyclic Graphs (PDAGs) [19] (see Methods) in which some undirected edges are present along with the directed ones as well as the $p$-values of all conditional independence tests. For a given PDAG detected by MRPC in each replicate and scenario, we utilise [27] to estimate the causal effects between the exposures and outcomes. Finally, Partition-DAG (ParDAG) [30] provides a solution to the structure learning problem once the summary-level statistics have been partitioned into two groups and the orientation of the edges from the exposures to the outcomes has been enforced. ParDAG computes the causal effects estimates under Lasso regularisation. It has not been combined with instrumental variable estimation and applied to genetic data to date. All methods use summary-level statistics as input after IVW. Finally, for each method and algorithmic implementation, details of the parameter settings are provided in Supplementary Information.

Regarding the evaluation criteria, we use a precision-recall curve (PRC) that shows the relationship between precision (*i.e.*, positive predictive value, on the $y$-axis) and recall (*i.e.*, sensitivity, on the $x$-axis) for every possible cut-off and it is not impacted by the over-representation of null effects. See Supplementary Information for a detailed discussion regarding how we implemented a fair comparison between the methods considered.

Finally, to evaluate the quality of the causal effects estimation, we calculate the sum of squared errors (SSE), defined as the sum of the squared differences between the estimated

and the simulated causal effect. In contrast to the evaluation of the recovery obtained by each method of the simulated dependencies within the exposures, the outcomes and between them, we do not report the SSE of the mediation parameters $\mathbf{\Gamma}_X$ and $\mathbf{\Gamma}_Y$ since they are considered nuisance parameters in the proposed model (see Supplementary Information).

## MrDAG more accurately detects unconfounded dependency relations within the exposures and the outcomes and between them

Figure 3 presents the results of MrDAG and alternative methods for one replicate of the simulated scenario UndG$_X$-DAG$_Y$ (Figures 3A-D) and DAG$_X$-DAG$_Y$ (Figures 3E-F) for a particular choice of the parameters $r_X = 0.6$ and $m_Y = 1$ used in the simulation study to control the average value of the mediation parameters $\mathbf{\Gamma}_X$ within the exposures and $\mathbf{\Gamma}_Y$ within the outcomes, and $\psi_X = 2$ and $\psi_Y = 1$ for the level of confounding on the exposures and the outcomes, respectively (see Methods).

The general performance of competing algorithms is already apparent from it. In scenario UndG$_X$-DAG$_Y$, if a causal effect is simulated from an exposure to an outcome and there are dependency relations from this outcome to other responses (Figure 3A), MR-BMA adds erroneously causal effects to all linked responses with severe FP inflation (Figure 3B, FPs between $\widehat{X}_{12}^*$ and $\widehat{Y}_3^*$, $\widehat{Y}_4^*$, $\widehat{Y}_5^*$ depicted in red). On the other hand, MR-BMA estimates neither the dependency pattern within $\boldsymbol{X}$, since (partial) correlation between summary-level exposures is assumed in the model [3] but not estimated, nor the dependencies within $\boldsymbol{Y}$ since MR-BMA considers one response at-a-time. MRPC infers correctly most of the dependencies within $\boldsymbol{X}$, but it does not have the power to detect all simulated causal effects $\boldsymbol{\Theta}$ at the specified type I error rate for the conditional independence test ($\alpha = 0.01$) with a few FNs (Figure 3C, FNs between $\widehat{X}_1^*$, $\widehat{X}_2^*$ and $\widehat{Y}_2^*$) and well as FPs within $\widehat{\boldsymbol{B}}_Y^*$ (FPs between $\widehat{Y}_2^*$, $\widehat{Y}_3^*$, $\widehat{Y}_4^*$, $\widehat{Y}_5^*$, where bidirectionally is erroneously detected). MrDAG performs better than alternative methods to detect both directed and bidirected edges with only one FP between $\widehat{X}_5^*$ and $\widehat{X}_{15}^*$ (Figure 3D).

Similar comments can be made for a particular replicate of scenario DAG$_X$-DAG$_Y$, although in this scenario the dependency patterns are more complex since a topological ordered DAG is simulated also within the outcomes (Figure 3E). MrDAG confirms its good performance except for the directionality of the dependency relations within $\boldsymbol{X}$, where bidirectional edges are found with a few FPs (Figure 3H, FPs between $\widehat{X}_1^*$ and $\widehat{X}_{12}^*$ and between $\widehat{X}_8^*$ and $\widehat{X}_9^*$).

Figure 4 generalises the results depicted in Figure 3, averaging the results over 25 replicates of the simulated scenarios UndG$_X$-DAG$_Y$ (Figures 4A-C) and DAG$_X$-DAG$_Y$ (Figures 4D-F) with the same parameters setting used in Figure 3. The results are
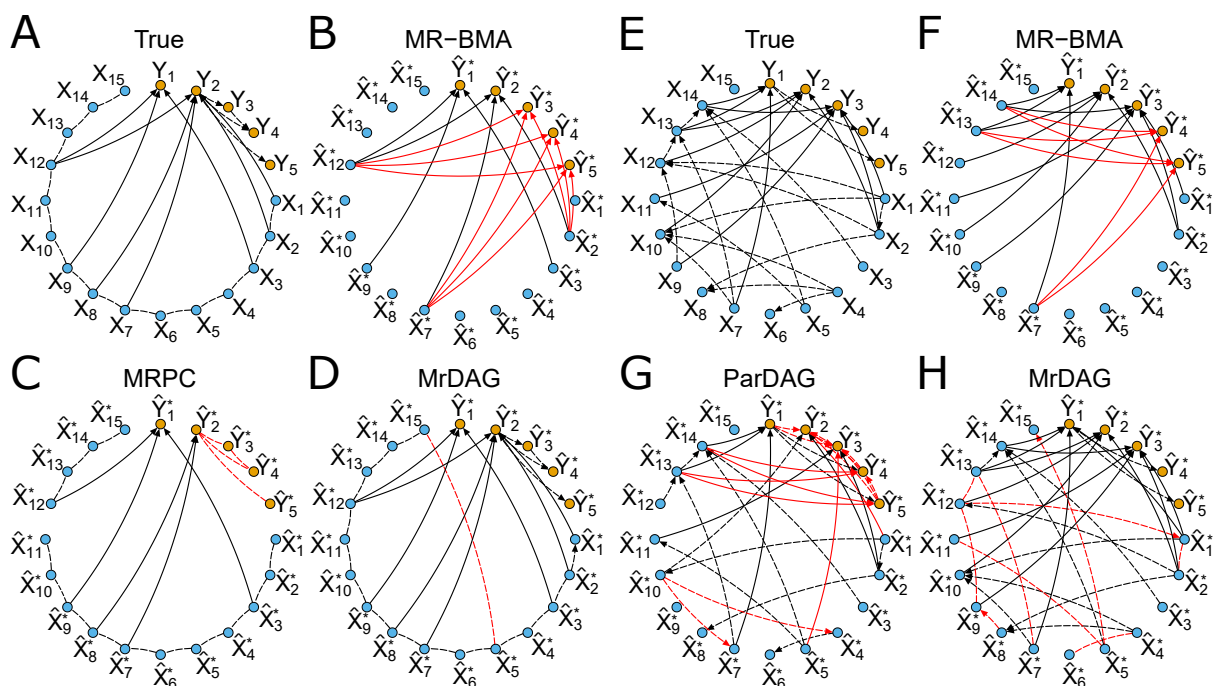
**Figure 3. Examples of unconfounded dependency structure simulated at the individual-level data and estimated by using summary-level statistics within the exposures, the outcomes and between them in two different scenarios.** In each panel, individual-level outcomes $Y = (Y_1, \ldots, Y_5)$ and exposures $X = (X_1, \ldots, X_{15})$ as well as genetically predicted outcomes $\widehat{Y}^* = (\widehat{Y}_1^*, \ldots, \widehat{Y}_5^*)$ and exposures $\widehat{X}^* = (\widehat{X}_1^*, \ldots, \widehat{X}_{15}^*)$ are represented with orange and blue nodes, respectively. Directed edges indicate dependency relations, while undirected edges denote partial correlation. Dashed lines depict the true (unconfounded by $U$) and estimated dependency structure within the exposures and the outcomes, while solid lines indicate true and estimated causal effects between them. Red colour denotes false positives, either falsely detected effects (regardless of the directionality) or wrong directionality of the edges. Besides the proposed model, alternative methods considered: Mendelian randomization with Bayesian Model Averaging (MR-BMA) [1], Mendelian randomization with PC algorithm (MRPC) [28], Partition-DAG (ParDAG) [30]. We report the results of MR-BMA and MrDAG without any threshold on the marginal posterior probability of inclusion (mPPI) and the posterior probability of edge inclusion (PPeI), respectively. MRPC Partially Directed Acyclic Graphs (PDAGs) are obtained by specifying the type I error rate for the conditional independence test at $\alpha = 0.01$. ParDAG results are the solutions of causal effects estimation with Lasso penalisation set at $\lambda = 0.9$ after partitioning the traits into two groups and enforcing a constraint on the orientation of the edges between the exposures and the outcomes. (**A**-**D**) Single replicate of the simulated scenario UndG$_X$-DAG$_Y$, where an undirected graph encodes the dependency pattern within $X$ and a DAG represents the dependency relations within $Y$ along with the simulated causal effects from the exposures to the outcomes, resulting in an overall partially oriented DAG. In this scenario, the strength of correlation between consecutive $X$ is set at $r_X = 0.6$, and then decreases exponentially for non-consecutive exposures, and the average level of the mediation parameters within $Y$ is set at $m_Y = 1$. (**E**-**H**) Single replicate of the simulated scenario DAG$_X$-DAG$_Y$, where two topologically ordered DAGs have been independently simulated within $X$ and $Y$ along with the simulated causal effects from the exposures to the responses, resulting in an overall fully-oriented DAG. In this scenario, the average level of mediation parameters for $X$ and $Y$ are set at $r_X = 0.6$ and $m_Y = 1$, respectively.

presented separately for the simulated dependency structures from the exposures to the outcomes (Figures 4A and D), within the exposures (Figures 4B and E) and within the outcomes (Figures 4C and D), respectively.

On average, MRPC and MrDAG have good performance in both simulated scenarios (Figures 4A and D). MRPC best results are obtained at a stringent type I error rate
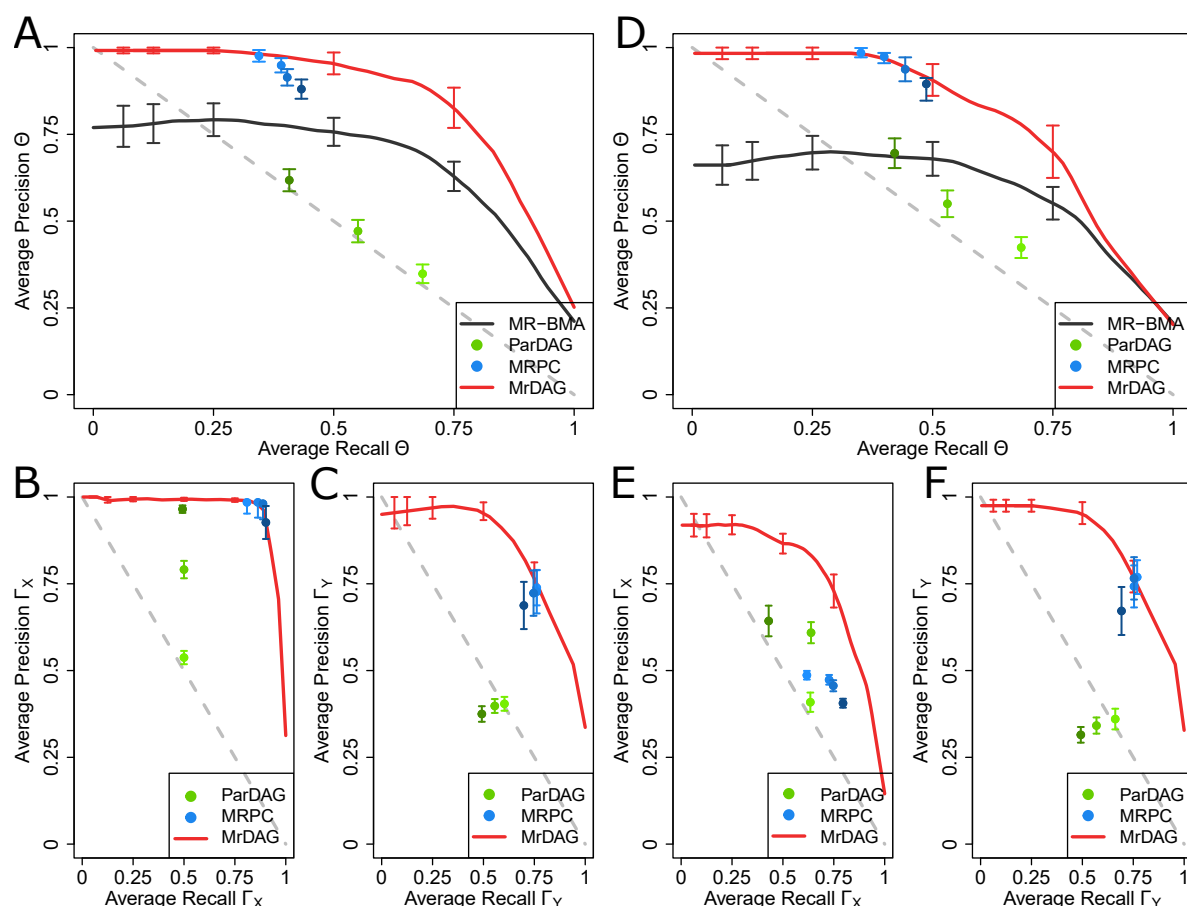
**Figure 4. Precision Recall Curves (PRCs) for all methods considered in the simulated scenarios UndG$_X$-DAG$_Y$ and DAG$_X$-DAG$_Y$** show recall (= sensitivity = TP/(TP + FN)) in the $x$-axis and precision (= positive predictive value = TP/(TP + FP)) in the $y$-axis with TP = True Positive, FN = False Negative and FN = False Positive averaged over 25 replicates in each scenario. In scenario UndG$_X$-DAG$_Y$ (**A-C**), the strength of correlation between consecutive $X$ is set at $r_X = 0.6$, and then it decreases exponentially for non-consecutive exposures, and the average level of the mediation parameters within $Y$ is set at $m_Y = 1$, while in scenario DAG$_X$-DAG$_Y$ (**D-F**), the average level of the mediation parameters within $X$ and $Y$ is set at $r_X = 0.6$ and $m_Y = 1$, respectively. For details, see Methods. In both scenarios, the results are presented separately for the simulated dependency structures from the exposures to the outcomes (A and D), within the exposures (B and E) and the outcomes (C and D), respectively. Vertical bars in each PRC, at specific recall levels $0.0625, 0.125, 0.25, 0.50$ and $0.75$, indicate standard error. For the MRPC algorithm, type I error rate for the conditional independence test is set at $\alpha = \{0.01, 0.05, 0.10, 0.20\}$ (from light- to dark-blue dots) and for the ParDAG algorithm we specify three different values for the Lasso penalisation $\lambda = \{0.5, 0.7, 0.9\}$ (from light- to dark-green dots). See Supplementary Information for details.

$\alpha = 0.01$ for the conditional independent tests (blue dots) although they are quite similar across different values of $\alpha$ and thus robust to this choice. However, it fails to detect the simulated dependency pattern within $X$ in scenario DAG$_X$-DAG$_Y$ (Figure 4B). The performance of MR-BMA can be only evaluated for the detection of the causal effects from the exposures to the outcomes (Figures 4A and D). As we noticed above, the large number of FPs degrades the results of this method which was not developed to deal with multiple related responses.

The performance of ParDAG is the worst among the methods considered for all types of designed relationships, slightly better within the exposures (Figures 4B and E) and

between the exposures and outcomes (Figures 4A and D) and worse within the outcomes (Figures 4C and F). Since ParDAG detects only directed edges, in Figure 4B, where partial correlation between exposures is simulated, the method has 50% recall rate. The results seem also quite different according to the penalty parameter $\lambda$.

MrDAG has a strong performance in both scenarios. In contrast to MR-BMA, in scenario $\mathrm{DAG}_X$-$\mathrm{DAG}_Y$ (Figures 4A and C) there is only a small reduction of the precision in the estimation of the dependency relations between the exposures and the outcomes, and within the latter, compared to the scenario $\mathrm{UndG}_X$-$\mathrm{DAG}_Y$ (Figures 4D and F).

The comments above can be extended to the scenarios where the relationships within outcomes are completely mediated ($\mathrm{UndG}_X$-$\mathrm{Med}_Y$ depicted in Supplementary Figures 2A-C, and $\mathrm{DAG}_X$-$\mathrm{Med}_Y$ shown in Supplementary Figures 2D-F). In these scenarios, the mediation within the outcomes is easier to detect (Supplementary Figures 2C and F) than a topologically ordered DAG simulated within $\boldsymbol{Y}$.

Supplementary Figure 3 shows the results of the AUCPR to detect the causal effects $\boldsymbol{\Theta}$ and the sensitivity of the methods to different specifications of $r_X$ and $m_Y$. MrDAG confirms to be uniformly the best method with stable AUCPR for any combination of $r_X$ and $m_Y$ with similar AUCPR when partial correlation or a topological ordered DAG is simulated within $\boldsymbol{X}$ (Supplementary Figures 3A and B). MR-BMA performs well, especially in the scenario $\mathrm{UndG}_X$-$\mathrm{Med}_Y$ (Supplementary Figure 3A) which is the scenario that is most compatible for this method as well as in scenario $\mathrm{DAG}_X$-$\mathrm{Med}_Y$ (Supplementary Figure 3C), where its performance slightly decreases. Both MRPC and ParDAG seem to be less precise at higher levels of $r_X$ irrespective of the simulated scenario, with ParDAG also influenced by the value of $m_Y$. Similarly, Supplementary Figures 4) and 5 show the sensitivity of the algorithms to detect the simulated patterns within $\boldsymbol{X}$ and within $\boldsymbol{Y}$ for different specifications of $r_X$ and $m_Y$.

## MrDAG improves the estimation of the causal effects over existing methods

Figure 5A shows the Sum of Squares Error (SSE) of the causal effects $\boldsymbol{\Theta}$ between the exposures and the outcomes for all methods considered in the simulated scenario $\mathrm{UndG}_X$-$\mathrm{DAG}_Y$ and Figure 5B for the simulated scenario $\mathrm{DAG}_X$-$\mathrm{DAG}_Y$ across 25 replicates in each scenario with the same parameter setting and implementation of algorithms described above. For MRPC and ParDAG algorithms, we only show the results obtained at type I error rate for the conditional independence test $\alpha = 0.01$ and Lasso penalisation $\lambda = 0.9$, respectively. These values provide the best results for the two algorithms as shown in Figure 4 and Supplementary Figure 2.

MrDAG has the lowest SSE mean and median (white dots and horizontal black line, respectively) in both scenarios. As expected, when a topological ordered DAG is simulated within the exposures (Figure 5B), the violin plot have a wider range, showing more variable
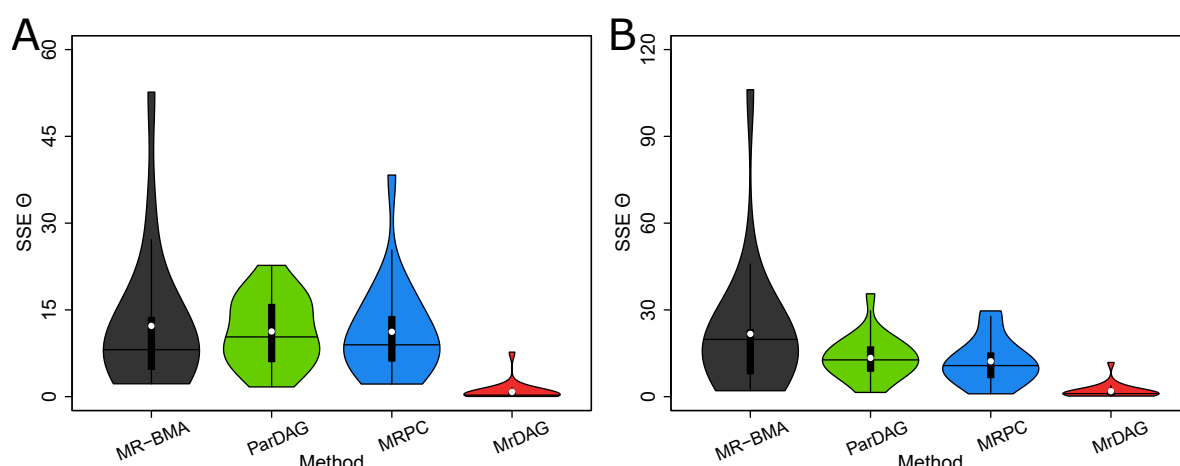
**Figure 5. Violin plots of the Sum of Squares Error (SSE) of the causal effects $\Theta$ between the exposures and the outcomes for all methods considered in the simulated scenarios $\text{UndG}_X$-$\text{DAG}_Y$ and $\text{DAG}_X$-$\text{DAG}_Y$ across 25 replicates in each scenario. (A)** In scenario $\text{UndG}_X$-$\text{DAG}_Y$, the strength of correlation between consecutive $X$ is set at $r_X = 0.6$, and then it decreases exponentially for non-consecutive exposures, and the average level of the mediation parameters within $Y$ is set at $m_Y = 1$. **(B)** In scenario $\text{DAG}_X$-$\text{DAG}_Y$, the average level of the mediation parameters within $X$ and $Y$ is set at $r_X = 0.6$ and $m_Y = 1$, respectively. For details, see Methods. In each violin plot, the vertical black thick line displays the interquartile range, the black horizontal line denotes the median and the white dot the mean. For MRPC and ParDAG algorithms, we only show the results obtained at type I error rate for the conditional independence test $\alpha = 0.01$ and Lasso penalisation $\lambda = 0.9$, respectively. These values provide the best results for the two algorithms as shown in Figure 4 and Supplementary Figure 2.

results, although the median is almost similar to the scenario with simulated partial correlation within $X$ (Figure 5A). Alternative methods have larger SSE.

Similar comments can be made for simulated scenarios $\text{UndG}_X$-$\text{Med}_Y$ (Supplementary Figure 6A) and $\text{DAG}_X$-$\text{Med}_Y$ (Supplementary Figure 6B), where a complete mediation is considered within the outcomes. MrDAG is confirmed as the best method.

We conclude this section by inspecting the sensitivity of the SSE of the causal effects between the exposures and the outcomes for different values of the average level of the mediation parameters $r_X$ and $m_Y$. The estimation of the causal effects displayed in Supplementary Figure 7 shows that both MR-BMA and MRPC depend on the combination of $r_X$ and $m_Y$ with similar performance when a complete mediation is simulated (Supplementary Figures 7A and C) (Supplementary Figures 7B and D). Compared to the other methods, MrDAG is not only the best, but it is rather insensitive to different levels of the mediation parameters within $X$ and $Y$.

# Real data application: The impact of lifestyle and behavioural traits on mental health

We apply MrDAG to investigate its ability to detect the effect of lifestyle and behavioural exposures on the risk of mental health phenotypes as well as potential forms of interventions for their prevention. As exposures, we chose seven lifestyle and behavioural traits that have previously been investigated for their effects on mental health, including educa-

tion (in years) (EDU), physical activity (PA), sleep duration (SP), alcohol consumption (ALC), lifetime smoking index (SM) and leisure screen time (LST). As outcomes, we select seven mental health phenotypes, including major depressive disorder (MDD), anorexia nervosa (AN), attention deficit hyperactivity disorder (ADHD), bipolar disorder (BD), autism spectrum disorder (ASD), schizophrenia (SCZ) and cognition (COG). See Supplementary Table 1 for the description of the summary-level statistics, the data sources, the number of IVs for each trait and Methods for the pre-processing steps. In a separate analysis, we also investigate the reverse direction, *i.e.*, whether the same mental health phenotypes have an impact on the group of lifestyle and behavioural traits by selecting IVs for the mental health phenotypes, see Methods for the respective pre-processing steps.

Figure 6 presents the results of MrDAG. In particular, Figures 6A and C show the estimated posterior probability of edge inclusion (PPeI) (12) after structure learning and Figures 6B and D the posterior causal effects (95% credible intervals (CI)) between the exposures and the outcomes. Results on PPeI (and the posterior causal effects) are not thresholded and sparsity is enforced by assigning a prior on the number of expected edges. We set it at $\pi^{\mathrm{edge}} = 0.16$, *i.e.*, we expect *a priori* one edge for each of the 13 traits, see Methods and Supplementary Information.

As shown in Figures 6C and D, there is one distinct exposure (LST) and two key shared exposures with important down-stream effects on mental health phenotypes, which are EDU and SM on which we focus our discussion. For each of them, we also describe how MrDAG can disentangle complex dependency relations within the exposures and the outcomes and detect (partial or complete) mediation which prevents spurious findings.

As could be expected due to its centrality in the global health agenda [31] and the high level of confounding of this phenotype with other genetically associated biological, behavioural and socioeconomic traits, genetically predicted EDU shows the most inter-exposure and exposure-outcome dependency relations (Figure 6C bottom part). Previous work has supported the broad mental health implications of education [32]. First, in keeping with previous findings [33, 34, 35, 36], our results show that EDU has a positive causal effect on COG, it is causally associated with an increased liability to ASD and BD as well as with a lower liability to ADHD. CIs show that the causal association with BD is markedly skewed to the right. In contrast, EDU has no effects on SP, the amount of ALC, or the liability to MDD [33], AN [37], or SCZ [36] (Figure 6D). Second, we investigate the detected dependency relations of EDU with other exposures that contribute to the reported causal associations. We find bidirectional relationships between genetically predicted EDU, PA and LST consistent with a large literature [33, 38]. Dependency relations have been also identified between EDU and SM [33, 39]. Supported by the existing literature, these results confirm the ability of MrDAG to disentangle complex relationships that exist between interrelated exposures.
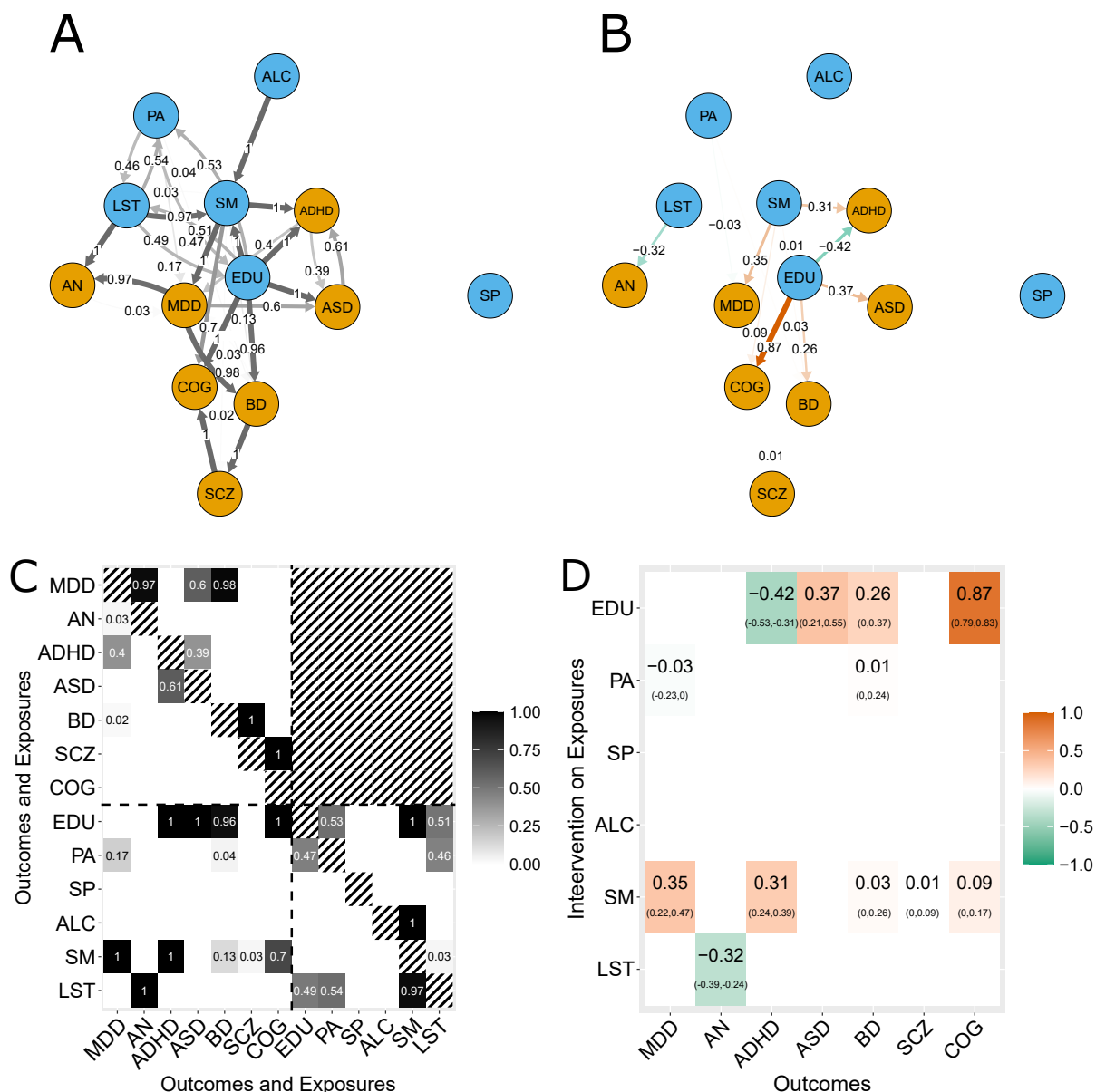
15

**Figure 6. Results of MrDAG algorithm regarding how lifestyle and behavioural exposures impact mental health outcomes.** (**A**) PDAG of the posterior probability of edge inclusion (PPeI) within the exposures (lifestyle and behavioural traits, blue nodes), the outcomes (mental health phenotypes, orange nodes) and between them. Undirected edges are represented as bidirectional edges, see, for instance, edges between PA (physical activity) and LST (leisure screen time) or ASD (autism spectrum disorder) and ADHD (attention deficit hyperactivity disorder). Neither reverse causation from the outcomes to the exposures nor feedback loops are allowed. (**B**) Posterior causal effects on the outcomes (orange nodes) under intervention on the exposures (blue nodes). Red and green edges indicate positive and negative posterior causal effects, respectively. (**C**) Posterior probability of edge inclusion (PPeI) for each combination of outcomes (mental health phenotypes) and exposures (lifestyle and behavioural traits). Horizontal and vertical dotted lines separate the exposures (bottom-right submatrix) from the outcomes (top-left submatrix). PPEIs between exposures and outcomes are depicted in the bottom-left submatrix. Neither reverse causation (top-right submatrix) nor feedback loops (main diagonal) are allowed (black-white strips). (**D**) Posterior causal effects (95% credible intervals) on the outcomes ($y$-axis) under intervention on the exposures ($x$-axis).

We find that SM is second only to EDU in its causal association with several outcomes. Specifically, SM associates with an increased liability to MDD and ADHD as previously reported [40, 41]. It is also associated with COG, BD and SCZ, although these causal

effects are small and CIs are skewed to the right. As discussed above, we also check the detected dependency relations of SM with other exposures. SM is related to PA as documented in epidemiological studies [42] and in standard MR analysis [43], the latter for objectively assessed average activity and number of cigarettes per day, respectively. Moreover, MrDAG appropriately identifies the relationship between ALC and SM, but not *vice versa*. In a recent MR publication [44], the opposite causal association is observed. However, in contrast to [44] who conceptualize SM with smoking initiation, we use a lifetime smoking index [40] which captures smoking duration, heaviness and cessation.

As important as the discussion of existing causal associations between the exposures and the outcomes, it is similarly insightful to discuss the absence of causal effects, especially those relationships that are reported in the literature or found by standard (one exposure and one outcome) MR models. For example, we do not replicate all previous evidence for positive causal effects of liability to SM on mental health phenotypes. Though we find a strong causal effect of SM on MDD [40], we do not find the same strong effect of SM on SCZ [40] as observed in observational studies [45, 46]. By looking at Figure 6C, this might be due to pleiotropic effects that have been identified by MrDAG within the mental health phenotypes. In line with prior findings, evidence from MrDAG supports dependency relations between genetic liability to MDD and AN, ASD and BD [47] as well as between genetic liability to BD and SCZ [48]. Lastly, in keeping with prior findings of possible bidirectional ASD-ADHD relationships [49], we observed genetic dependency relations between ASD and ADHD, and *vice versa*. These results suggest that the genetic effects of SM on SCZ can be mediated by pleiotropic effects within the responses. By considering the results above, we hypothesise that the SM to SCZ relationship is partly mediated first by MDD and then by BD. Moreover, there is another path that goes from the genetically predicted level of SM to SCZ through a positive weak causal association identified by MrDAG between SM and BD [50]. Both genetic paths are illustrated in Figure 6A. Conditionally on these relationships that are not considered in standard MR or MVMR, MrDAG does not detect a strong causal effect between SM and SCZ.

We further note that the causal effect of SM on ADHD is both direct and indirect, the latter mediated first by MDD and then by ASD. Thus, our analysis pinpoints the important role of MDD which partly or entirely accounts for many causal pathways within mental health phenotypes and their causal exposures. This might be due to the potentially high levels of confounding and non-specific genetic associations present in the original MDD GWAS [51, 52] as well as the high levels of symptom-level and therefore diagnostic overlap between MDD and all other psychiatric disorders [53]. Nonetheless, the implications of our results, assuming the validity of all GWAS findings, are that prevention and/or therapeutic intervention on MDD [54] can have a cascade of important effects for the prevention of several mental health phenotypes.

To investigate this hypothesis, Supplementary Figures 10A and B show the results

of MrDAG when MDD is removed from the list of outcomes. Regarding the causal association between SM and ADHD, it is still present with the same strength and similar CI depicted in Figure 6D, suggesting that the indirect effect mediated first by MDD and then by ASD is negligible. Supplementary Figure 10B also shows that, after removing MDD, the genetically predicted SM is positively associated with SCZ as reported in the literature. Combined with our main findings, this result indicates that the absence of a link between SM and SCZ in the MrDAG model is likely due to the mediation of MDD and BD.

The risk of detecting spurious shared causal effects is very high when a standard MR method is used separately on each trait as well as when multiple exposures are considered for each outcome [1]. This problem has been highlighted in the simulation study and visually presented in Figures 3B and F. In Supplementary Table 2 we show the results MR-BMA algorithm when applied to the same data set. We notice an overestimation of the causal effects since MR-BMA tries to ascribe the whole effects to the exposures and, as expected from the simulation study, it also detects many more associations than MrDAG.

We conclude the analysis by assessing the validity of the results obtained by MrDAG. We divide this internal check into sensitivity to hyper-prior specification and robustness of structure learning. Regarding the first point, Supplementary Figure 11 show that the posterior causal effects as well as the 95% CIs for different values of the *a priori* probability of edge inclusion are not influenced by this choice. For the second internal check, we bootstrap MrDAG repeatedly on the data [55] (see Supplementary Information). In Supplementary Figures 12 we present the bootstrap frequency of edge inclusion for each permitted combination of exposures and outcomes and the scatterplot of the posterior probability of edge inclusion (PPeI) against the bootstrap frequency of edge inclusion. The results show that there is a satisfactory agreement between a single run of the algorithm and the bootstrap results for the causal associations. Extended results are presented in Supplementary Information.

For completeness, we have also tested reverse causation by selecting genetic variants to be associated with the mental health phenotypes. Figure 7 and Supplementary Figure 9 show the results of the analysis to detect the impact of mental health phenotypes on lifestyle and behavioural traits, where, besides the positive causal effect of genetically predicted COG on EDU [56], smoking is causally affected by the genetic liability to MDD [40] and ADHD, the latter well-documented in epidemiological studies [57] and recently confirmed in a randomised clinical trial of smoking cessation [58].
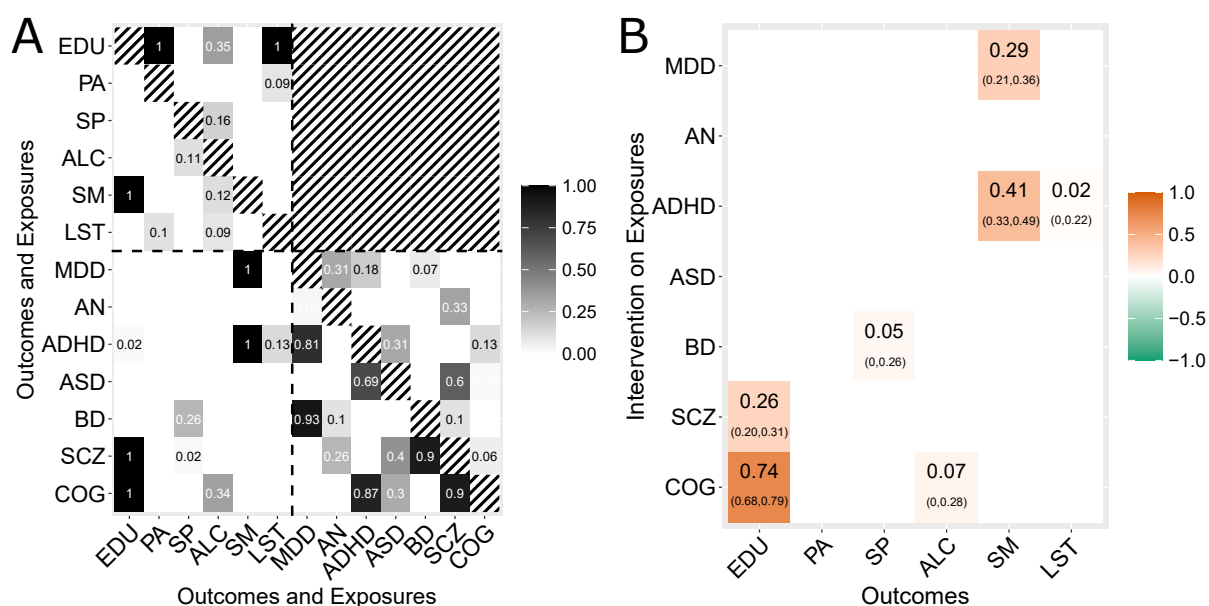
**Figure 7. Results of MrDAG regarding how liability to mental health phenotypes affects lifestyle and behavioural traits.** (**A**) Posterior probability of edge inclusion (PPeI) for each combination of outcomes (lifestyle and behavioural traits) and exposures (mental health phenotypes). Horizontal and vertical dotted lines separate the exposures (bottom-right submatrix) from the outcomes (top-left submatrix). PPEIs between exposures and outcomes are depicted in the bottom-left submatrix. Neither reverse causation (top-right submatrix) nor feedback loops (main diagonal) are allowed (black-white strips). (**B**) Posterior causal effects (95% credible intervals) on the outcomes (*y*-axis) under intervention on the exposures (*x*-axis).

# Discussion

Here, we have introduced MrDAG, the first Bayesian causal graphical MR model for multivariable and multiresponse that can detect dependency patterns within the exposures as well as within the outcomes thus allowing for a more precise estimation of the causal effects from the exposures to the outcomes. We showcased the advantage of the proposed method in a comprehensive simulation study and its utility in detecting how lifestyle and behavioural traits interact to cause mental health phenotypes, and *vice versa*. In the real data application, we highlighted how MrDAG can recover more information on the genetic paths that link exposures to outcomes compared to existing MR methods that ignore these dependency relations. Specifically, we highlighted education and smoking as key effective points of intervention given their distinct downstream effects on multiple mental health phenotypes.

These insights are possible since three methodological advances are considered in MrDAG. First, in structure learning, the hypothesis of no unobserved confounding is a fundamental underlying assumption. This assumption, known as causal sufficiency, is difficult to justify in real data applications and its violation produces biased results. By using IVs within the MR paradigm, we bypass the need to remove the effects of the unmeasured confounder from the individual-level data [21]. Instead, we solve this problem by employing genetically predicted exposures and outcomes which depend only on the

genetic variants chosen as IVs. Genetically predicted exposures are key in the derivation of the two-stage least square causal effect estimator [25], but in MrDAG we have extended it to include genetically predicted outcomes. On both predicted traits, we perform DAG exploration to learn the unconfounded dependency relations that exist within the exposures, the outcomes and between them. Our second contribution is the estimation of causal effects under intervention on the exposures conditionally on a given DAG. We showed that they can be estimated based on Pearl's interventional calculus [7]. Moreover, differently from [59] and its application in the MRPC algorithm [60], the estimation of the causal effects is averaged over the visited graphical models [61], thus taking into account the uncertainty regarding the graphs that best portray the dependency structure in a given data set. Third, MrDAG allows the possibility of including domain-knowledge relations between the traits. In the designed MrDAG model, constraints between the exposures and the outcomes descend directly from the MR paradigm. Our Bayesian implementation of structure learning under restrictions offers clear advantages over alternative methods [30]. Although not discussed here, other restrictions can be straightforwardly included, for instance, known relations regarding disease progression or time-dependent outcomes, *e.g.*, smoking initiation and cessation [62].

In the real data application, while the use of existing summary-level statistics of genome-wide association studies facilitates the integration of diverse phenotypes measured in different cohorts, we are also limited by the biases suffered by the initial genome-wide association studies. Specifically, studies on mental health rely on the presence of a clinical diagnosis. Consequently, it is not truly the genetic liability of the disease itself as much as it is the probability of having access to diagnoses or treatment. Our findings on the relationship between higher genetically predicted educational attainment (EDU) and increased ASD and BD, but decreased ADHD risk provide an example of such bias. In these analyses, the predicted number of school years completed is unlikely to be causally implicated in the development of ASD traits. While the typical age of onset of ASD precedes the start of formal education (therefore unlikely to be caused by it), ASD-related traits are more likely to be recognized and referred, particularly in those who are undiagnosed or untreated, when individuals are within a schooling system where standardized testing and progress reports by peer comparison are performed. Moreover, current GWAS consider one trait or disease at-a-time and do not consider to what extent cases are comorbid with other diseases. Future GWAS on co-morbidity [63] may provide more fine-grained genetic associations allowing to disentangle some of these relationships. Alternatively, novel causal inference methodology designed for individual-level data in combination with large-scale biobank or cohort studies with genotype data could be used to triangulate evidence.

In conclusion, MrDAG represents an important step forward in how we can learn complex relationships among phenotypic traits and uncover causal pathways using ge-

netic data. It provides analysts with the opportunity to derive a more comprehensive picture of causal mechanisms between complex phenotypes. The real data application is an example of the proposed holistic approach, where we leverage MrDAG and large-scale genome-wide association data to offer novel mechanistic insight into the causal behavioural determinants of mental health phenotypes to delineate between their overlapping pathophysiology and phenotypic presentation, toward translational progress in the field of mental health. Moving forward, MrDAG is ideally placed for the analysis of common causal exposures for multimorbid health conditions. This research into multimorbidity has been facilitated by the advent of large-scale biobanks being linked and followed up using electronic health records and routinely collected health care data. Using genotype data as genetic anchors offers a principled way for causal inference. MrDAG provides an addition to existing toolkits to map shared and distinct causes of disease, to understand trajectories, and to draw causal paths that link diseases.

# Methods

In the following, we denote with capital letters the random variables $Y$, $X$, $G$ and $U$ for the observed outcome, exposure, instrumental variable and unmeasured confounder, respectively, and with small letters $y$, $x$, $g$ and $u$ their corresponding observations. Multivariate random variables and corresponding observations are presented in bold. A marginal element of a vector of random variables is specified by a suitable subscript index, $e.g.$, $Y_k$, $k \in K = \{1, \ldots, q\}$, $X_j$, $j \in J = \{1, \ldots, p\}$, and $G_i$, $i \in I = \{1, \ldots, n\}$. $\boldsymbol{Y}_{\backslash k}$ and $\boldsymbol{Y}_{\backslash j}$ consists of all the outcomes and exposures except those that are related to the $k$th response and $j$th exposure, respectively. Finally, vectors understood as columns vectors and matrices are indicated in bold, the latter also in capital letters.

We indicate with $\beta_{i,j}^X$ and $\beta_{i,k}^Y$ the effect of the genetic variant $i \in I$ on the exposure $j \in J$ and outcome $k \in K$, respectively, with $\boldsymbol{\beta}_{X_j}$ and $\boldsymbol{\beta}_{Y_j}$ the $n$-dimensional vector of genetic effects on the $j$th exposure and $k$th outcome, respectively, and, finally, with $\boldsymbol{B}_X$ and $\boldsymbol{B}_Y$ the $(n \times p)$- and $(n \times q)$-dimensional matrices of the genetic effects on all exposures and outcomes. $\theta_{j,k}$ denotes the causal parameter of interest, $i.e.$, the causal effect of $X_j$ on $Y_k$, and $\gamma_{h,j}^X$ and $\gamma_{h,k}^Y$ the mediation effect of $X_h$ on $X_j$, $h \neq j$ and $Y_j$ on $Y_k$, $h \neq k$, respectively. $\boldsymbol{\Theta}$, $\boldsymbol{\Gamma}_X$ and $\boldsymbol{\Gamma}_Y$ indicate the corresponding $(p \times q)$-, $(p \times p)$- and $(q \times q)$-dimensional matrices of the causal parameters of interest ($\boldsymbol{\Theta}$) and mediation parameters ($\boldsymbol{\Gamma}_X$ and $\boldsymbol{\Gamma}_Y$). The symbol "$\widehat{\phantom{x}}$" denotes the estimator of a parameter or its estimated value and "$*$" an IVW parameter.

Finally, let $\mathcal{D} = (V, E)$ be a Directed Acyclic Graph (DAG), where $V$ denote a set of vertices (nodes) and $E = V \times V$ a set of directed edges, $i.e.$, if $(z, v) \in E$, then $(z, v) \notin E$. For a given DAG $\mathcal{D}$, if $z \rightarrow v$, then $z$ is a parent of $v$ and, conversely, $v$ is a child of $z$. Moreover, if $z \rightarrow \ldots \rightarrow v$, then $z$ is an ancestor of $v$ and $v$ is a descendant of $z$. We

denote the parent set of $v$ in $\mathcal{D}$ as $\mathrm{pa}_{\mathcal{D}}(v)$ and $v \cup \mathrm{pa}_{\mathcal{D}}(v) = \mathrm{fa}_{\mathcal{D}}(v)$ the family of $v$. Unless otherwise stated, for ease of notation, we remove the subscript $\mathcal{D}$.

In [5, 20, 64] key results regarding standard Mendelian randomization (single exposure with single instrumental variable and single outcome) are presented. Here, we use them to show that MrDAG is an extension of standard MR when (i) multiple exposures and outcomes are considered and (ii) the underlying dependency relations within and between them are not known (latent) and need to be estimated from the data. Technical details are provided in Supplementary Information.

## Multi-exposure and multi-outcome core conditions for instrumental variables

Let $\boldsymbol{Y}$, $\boldsymbol{X}$ and $\boldsymbol{G}$ be the $q$-, $p$- and $n$-dimensional vector of the outcomes, exposures and instruments (genotypes) random variables, respectively.

Let's assume the following "multivariate core conditions" (MCC) for valid instrumental variables (IVs) which are the extensions of the core conditions that $\boldsymbol{G}$ has to satisfy in standard Mendelian randomization (MR) [5]:

(IV1) $G_i \perp\!\!\!\perp U$, $\forall i \in I$, *i.e.*, $G_i$ must be independent of $U$;

(IV2) $G_i \not\perp\!\!\!\perp X_j \mid \boldsymbol{X}_{\backslash j}$, $\forall i \in I$ and $\forall j \in J$, *i.e.*, $G_i$ must *not* be independent of $X_j$ conditionally on $\boldsymbol{X}_{\backslash j}$;

(IV3) $G_i \perp\!\!\!\perp Y_k \mid (\boldsymbol{X}, U)$, $\forall i \in I$ and $\forall k \in K$, *i.e.*, $G_i$ must be independent of $Y_k$ conditionally on $\boldsymbol{X}$ and $U$.

The first multi-exposure and multi-outcome core condition (MCC) for instrumental variables is similar to the first CC in standard MR [5]. The second MCC imposes that $G_i$ should be associated with $X_j$ conditionally on the other exposures. The third MCC establishes that the instrumental variables and outcomes are conditionally independent given the exposures and the unmeasured confounder.

From the DAG $\mathcal{D}$ involving $\boldsymbol{Y}$, $\boldsymbol{X}$, $\boldsymbol{G}$ and $U$ that satisfies the MCC, the corresponding Markov properties say that $G_i \perp\!\!\!\perp U$, $\forall i \in I$, since $G_i$ is not a descendant of $U$ and *vice versa* and $G_i \not\perp\!\!\!\perp X_j \mid \boldsymbol{X}_{\mathrm{pa}(j)}$, $\forall i \in I$ and $\forall j \in J$, because $X_j$ is a descendant of $G_i$. The Markov property for the third MCC is $G_i \perp\!\!\!\perp Y_k \mid (\boldsymbol{Y}_{\mathrm{pa}(k)}, \boldsymbol{X}_{\mathrm{pa}(k)}, U)$, $\forall i \in I$ and $\forall k \in K$, since $G_i$ is a non-descendant of $Y_k$ and $(\boldsymbol{Y}_{\mathrm{pa}(k)}, \boldsymbol{X}_{\mathrm{pa}(k)}, U)$ are the parents of $Y_k$.

## Interventional distributions and causal effects estimation

The conditional dependencies associated with the multi-exposure and multi-outcome DAG $\mathcal{D}$ lead to the following factorisation of the joint density of all random variables considered

$$f(\boldsymbol{y}, \boldsymbol{x}, \boldsymbol{g}, u) = \prod_{k \in K} f(y_k \mid \boldsymbol{y}_{\mathrm{pa}(k)}, \boldsymbol{x}_{\mathrm{pa}(k)}, u) \prod_{j \in J} f(x_j \mid \boldsymbol{x}_{\mathrm{pa}(j)}, \boldsymbol{g}, u) f(\boldsymbol{g}) f(u)$$

which is known as pre-intervention distribution and it is assumed to be faithful to the DAG [29], *i.e.*, there are no conditional dependence relationships between the variables in the model that do not follow directly from the Markov properties.

The post-intervention distribution under intervention on the $h$th exposure sets to take the value $\widetilde{x}_h$ is obtained by the truncated factorisation [7]

$$f(\boldsymbol{y}, \boldsymbol{x}_{\backslash h}, \boldsymbol{g}, u \mid \mathrm{do}(X_h = \widetilde{x}_h)) = \prod_{k \in K} f(y_k \mid \widetilde{x}_h, \boldsymbol{y}_{\mathrm{pa}(k)}, \boldsymbol{x}_{\mathrm{pa}(k)}, u) \\ \prod_{j \in J \backslash \{h\}} f(x_j \mid \boldsymbol{x}_{\mathrm{pa}(j)}, \boldsymbol{g}, u) \mathbb{I}_{x_h}(\widetilde{x}_h) f(\boldsymbol{g}) f(u), \tag{1}$$

where $\mathbb{I}_x(\widetilde{x})$ is the indicator function which is equal to one if $x_h = \widetilde{x}_h$ and zero otherwise. Graphically, the directed edges to $X_h$ from its parents in $\boldsymbol{X}$, $\boldsymbol{G}$ and $U$ are removed.

A post-intervention distribution under intervention on the $h$the exposure is obtained from (1) by marginalising all variables but the selected outcome and the exposure on which an intervention is carried out

$$f(y_k \mid \mathrm{do}(X_h = \widetilde{x}_h)) = \int f(\boldsymbol{y}, \boldsymbol{x}_{\backslash h}, \boldsymbol{g}, u \mid \mathrm{do}(X_h = \widetilde{x}_h)) \, \mathrm{d}\boldsymbol{y}_{\backslash k} \, \mathrm{d}\boldsymbol{x}_{\backslash h} \, \mathrm{d}\boldsymbol{g} \, \mathrm{d}u \\ = \int f(y_k \mid \widetilde{x}_h, \boldsymbol{x}_{\mathrm{pa}(h)}, u) f(\boldsymbol{x}_{\mathrm{pa}(h)}, u) \mathbb{I}_{x_h}(\widetilde{x}_h) \, \mathrm{d}\boldsymbol{x}_{\mathrm{pa}(h)} \, \mathrm{d}u. \tag{2}$$

This result is derived from [7] and it follows directly from the Markov properties of the DAG. It establishes that the parents of the variable on which an intervention is carried out are the only variables that need to be measured to estimate the causal effect on an outcome [65].

The post-intervention distribution (2) can be summarised by taking the expectation and defining the causal effect of an intervention [59] as

$$\theta_{h,k} = \frac{\partial}{\partial x_h} \mathbb{E}(Y_k \mid \mathrm{do}(X_h = x_h)) \bigg|_{x_h = \widetilde{x}_h}, \quad h \in J, k \in K.$$

In Supplementary Information, we show the identifiability of the causal effect (Supplementary Proposition 2) and the derivation of its estimand in multiple exposures and multiple outcomes MR framework (Supplementary Proposition 3). We also show the consistency of the effects of the regressions of each outcome and exposure on $\boldsymbol{G}$ (Supplementary Proposition 1), *i.e.*, the estimated genetic effects on the outcomes and exposures contain all information regarding the causal parameters of interest and the mediation

23

parameters within the exposures and the outcomes.

Here, for a given DAG $\mathcal{D}$, we report the IVW estimator of the causal effect of the intervention in $X_h$ on $Y_k$

$$\widehat{\theta}_{h,k} = [(\widehat{\boldsymbol{B}}_{X_{\mathrm{fa}(h)}}^{*\top} \widehat{\boldsymbol{B}}_{X_{\mathrm{fa}(h)}}^*)^{-1} \widehat{\boldsymbol{B}}_{X_{\mathrm{fa}(h)}}^{*\top} \widehat{\boldsymbol{\beta}}_{Y_k}^*]_1, \tag{3}$$

where the subscript indicates the first element of the solution of the linear least squares (LLS) regression since $\mathrm{fa}(v) = v \cup \mathrm{pa}(v)$, $\boldsymbol{X}_{\mathrm{fa}(h)}$ denotes the exposures that are the family of the exposure $X_h$ under intervention, $\widehat{\boldsymbol{B}}_{X_{\mathrm{fa}(h)}}^*$ are the IVW estimated coefficients of the regressions of each exposure in $\boldsymbol{X}_{\mathrm{fa}(h)}$ on $\boldsymbol{G}$ and $\widehat{\boldsymbol{\beta}}_{Y_k}^*$ is the IVW estimated coefficient of a regression of $Y_k$ on $\boldsymbol{G}$. (3) resembles the standard IVW estimator of the causal effect that approximates the estimate that would have been obtained if individual-level data were available [3]. However, in contrast to general proposed solutions in MVMR, in (3) the set of regressors is with regard to the family of the exposure under intervention.

## Dependency structure under the effect of unmeasured confounders

To estimate (3), structure learning of the graphical models needs to be performed to detect the parents $\boldsymbol{X}_{\mathrm{pa}(h)}$ of the exposure $X_h$ under intervention. However, structure learning assumes causal sufficiency [21], *i.e.*, it requires that there are no hidden (or latent) variables that are common causes of two or more traits. Instead, here we explicitly assume that an unmeasured confounder $U$ acts on both outcomes and exposures.

Links between the genetic correlation and MR causal effect estimate have been already discussed in [23]. Here, we provide further connections with genetic covariance [24] which is key to show that, by working with summary-level statistics, it is possible to recover the dependency structure between the corresponding traits in the original (individual-level) data unconfounded by $U$.

Let's assume that the genetic effect on a phenotypic trait is linear and consider two traits

$$Y_k = \boldsymbol{G}^\top \boldsymbol{\beta}_{Y_k} + \psi_Y U + \epsilon_{Y_k}, \quad k \in K,$$
$$X_j = \boldsymbol{G}^\top \boldsymbol{\beta}_{X_j} + \psi_X U + \epsilon_{X_j}, \quad j \in J,$$

where $\boldsymbol{G}$ is a set of genetic variants, either spanning the whole genome, or region(s)-specific or selected to be associated with a trait, $\boldsymbol{\beta}_{Y_k}$ and $\boldsymbol{\beta}_{Y_k}$ are the genetic effects, $U$ is an unmeasured confounder that affects both traits with $\psi_Y$ and $\psi_X$ the effects sizes and $\epsilon_{Y_k}$ and $\epsilon_{X_k}$ are white noises which can be interpreted as environmental effects. We assume that $\boldsymbol{G} \perp\!\!\!\perp U$ and, similarly, $\boldsymbol{G} \perp\!\!\!\perp \epsilon_{Y_k}$ and $\boldsymbol{G} \perp\!\!\!\perp \epsilon_{X_j}$. Finally, we assume that $U \perp\!\!\!\perp \epsilon_{Y_k}$ and $U \perp\!\!\!\perp \epsilon_{X_j}$, *i.e.*, the unmeasured confounder $U$ exerts its effect on both traits and it is distinct from other environmental factors. Under this model, the phenotypic

24

covariance is

$$
\begin{aligned}
\mathrm{Cov}(Y_k, X_j) &= \mathrm{Cov}(\boldsymbol{G}^\top \boldsymbol{\beta}_{Y_k} + \psi_Y U + \epsilon_{Y_k}, \boldsymbol{G}^\top \boldsymbol{\beta}_{X_j} + \psi_X U + \epsilon_{X_j}) \\
&= \boldsymbol{\beta}_{Y_k}^\top \mathbb{V}(\boldsymbol{G})\boldsymbol{\beta}_{X_j} + \psi_Y \psi_X \mathbb{V}(U) + \mathrm{Cov}(\epsilon_{Y_k}, \epsilon_{X_j}).
\end{aligned}
\tag{4}
$$

The phenotypic covariance can be decomposed into $c_g(Y_k, X_k) = \mathrm{Cov}(\boldsymbol{G}^\top \boldsymbol{\beta}_{Y_k}, \boldsymbol{G}^\top \boldsymbol{\beta}_{X_j}) = \boldsymbol{\beta}_{Y_k}^\top \mathbb{V}(\boldsymbol{G})\boldsymbol{\beta}_{X_j}$, the genetic covariance between the two traits, *i.e.*, the covariance between the genetic components of the two traits, $\boldsymbol{G}^\top \boldsymbol{\beta}_{Y_k}$ and $\boldsymbol{G}^\top \boldsymbol{\beta}_{X_j}$, and the environmental covariance, *i.e.*, the covariance between the environmental effects of two traits that we have split into the effect of the unmeasured confounder, $c_u(Y_k, X_k) = \psi_Y \psi_X \mathbb{V}(U)$, and other environmental factors, $c_e(Y_k, X_k) = \mathrm{Cov}(\epsilon_{Y_k}, \epsilon_{X_j})$. If the environmental factors are trait-specific since $U$ includes all common confounding factors, $c_e(Y_k, X_k) = 0$ and (4) shows that an estimand of the covariance between two traits unconfounded by $U$ is $c_g$. From an MR perspective, by using MCC with $\boldsymbol{G}$ a set of IVs, in Supplementary Proposition 5 we show that $\mathrm{Cov}(Y_k, X_h \mid \boldsymbol{G} = \boldsymbol{g})$ is unconfounded by $U$.

Assuming that the individuals for the two phenotypic traits are drawn from the same population with LD matrix between the genetic variants $\boldsymbol{V} = \boldsymbol{G}^\top \boldsymbol{G}$, the sampling distribution of the genetic effects are $N_{Y_k}^{1/2}(\widehat{\boldsymbol{\beta}}_{Y_k} - \boldsymbol{\beta}_{Y_k}) \xrightarrow{d} \mathrm{N}_n(\boldsymbol{0}, \sigma_{Y_k}^2 \boldsymbol{V}^{-1})$ and $N_{X_j}^{1/2}(\widehat{\boldsymbol{\beta}}_{X_j} - \boldsymbol{\beta}_{X_j}) \xrightarrow{d} \mathrm{N}_n(\boldsymbol{0}, \sigma_{X_j}^2 \boldsymbol{V}^{-1})$, where "$d$" denotes convergence in distribution. Under infinite sample sizes, $\widehat{\boldsymbol{\beta}}_{Y_k} \xrightarrow{p} \boldsymbol{\beta}_{Y_k}$ and $\widehat{\boldsymbol{\beta}}_{X_j} \xrightarrow{p} \boldsymbol{\beta}_{X_j}$, where "$p$" denotes convergence in probability, and an estimator of the genetic covariance between the two traits is

$$
\widehat{c}_g(Y_k, X_j) = \widehat{\boldsymbol{\beta}}_{Y_k}^\top \boldsymbol{V} \widehat{\boldsymbol{\beta}}_{X_j}.
$$

In the finite sample sizes case, the estimates of $\boldsymbol{\beta}_{Y_k}$ and $\boldsymbol{\beta}_{X_j}$ are noised and $\widehat{c}_g(Y_k, X_j)$ is biased [24]

$$
\mathbb{E}(\widehat{c}_g(Y_k, X_j)) = \boldsymbol{\beta}_{Y_k}^\top \boldsymbol{V} \boldsymbol{\beta}_{X_j} + \frac{N_o}{N_{Y_k} N_{X_j}} c_u(Y_k, X_j),
\tag{5}
$$

where $N_o$ is the sample size overlap between the two traits. However, even in the scenario of complete overlap, the bias in (5) is negligible if the sample sizes of the two traits are large, as it usually happens in modern GWAS.

The same considerations can made for all phenotypic traits under investigation to reconstruct their joint genetic covariance unconfounded by $U$

$$
\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{YY} & \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XY} & \boldsymbol{\Sigma}_{XX} \end{bmatrix} = \begin{bmatrix} \boldsymbol{B}_Y^\top \boldsymbol{V} \boldsymbol{B}_Y & \boldsymbol{B}_Y^\top \boldsymbol{V} \boldsymbol{B}_X \\ \boldsymbol{B}_X^\top \boldsymbol{V} \boldsymbol{B}_Y & \boldsymbol{B}_X^\top \boldsymbol{V} \boldsymbol{B}_X \end{bmatrix},
\tag{6}
$$

where $\boldsymbol{\Sigma}_{XX}$, $\boldsymbol{\Sigma}_{YY}$ and $\boldsymbol{\Sigma}_{XY}$ are the genetic covariances within the exposures, the outcomes and between them and $\boldsymbol{B}_Y$ and $\boldsymbol{B}_X$ are the coefficients of the regressions of the outcomes and the exposures on $\boldsymbol{G}$, respectively.

## MrDAG model

Assuming that the individuals for two phenotypic traits $Y_k$ and $X_j$ are drawn from the same population with LD matrix $\boldsymbol{V}$, we have $N_{Y_k}^{1/2}(\boldsymbol{g}^\top\widehat{\boldsymbol{\beta}}_{Y_k} - \boldsymbol{g}^\top\boldsymbol{\beta}_{Y_k}) \overset{d}{\to} \mathrm{N}(0, \sigma_{Y_k}^2)$ and $N_{X_j}^{1/2}(\boldsymbol{g}^\top\widehat{\boldsymbol{\beta}}_{X_j} - \boldsymbol{g}^\top\boldsymbol{\beta}_{X_j}) \overset{d}{\to} \mathrm{N}(0, \sigma_{X_j}^2)$, where $\boldsymbol{g}$ are the observed IVs, $\boldsymbol{g}^\top\widehat{\boldsymbol{\beta}}_{Y_k}$ and $\boldsymbol{g}^\top\widehat{\boldsymbol{\beta}}_{X_j}$ are the $k$th and the $j$th genetically predicted values of the outcome and exposure, *i.e.*, $\widehat{Y}_k$ and $\widehat{X}_k$, respectively.

The joint distribution of all genetically predicted values of the outcomes and exposures based on the IVs is

$$[\boldsymbol{g}^\top\widehat{\boldsymbol{B}}_Y \boldsymbol{g}^\top\widehat{\boldsymbol{B}}_X]^\top \sim \mathrm{N}_{q+p}([\boldsymbol{g}^\top\boldsymbol{B}_Y \boldsymbol{g}^\top\boldsymbol{B}_X]^\top, \boldsymbol{\Sigma}),$$

*i.e.*, for large sample sizes they are normally distributed with mean $[\boldsymbol{g}^\top\boldsymbol{B}_Y \boldsymbol{g}^\top\boldsymbol{B}_X]^\top$ and covariance matrix $\boldsymbol{\Sigma} \in \mathcal{C}_\mathcal{D}$, the space of the symmetric positive definite covariance matrices Markov with respect to the DAG $\mathcal{D}$.

If we assume that IVW is performed on the estimated regression coefficients and IVs are independent after pruning or clumping, *i.e.*, $\boldsymbol{V} = \boldsymbol{I}_n$, the MrDAG model becomes

$$[\boldsymbol{g}^\top\widehat{\boldsymbol{B}}_Y^* \boldsymbol{g}^\top\widehat{\boldsymbol{B}}_X^*]^\top \sim \mathrm{N}_{q+p}([\boldsymbol{g}^\top\boldsymbol{B}_Y^* \boldsymbol{g}^\top\boldsymbol{B}_X^*]^\top, \boldsymbol{\Sigma}^*), \tag{7}$$

where $[\boldsymbol{g}^\top\widehat{\boldsymbol{B}}_Y^* \boldsymbol{g}^\top\widehat{\boldsymbol{B}}_X^*]^\top = [\boldsymbol{g}^\top\overline{\sigma}_Y^{-1}\widehat{\boldsymbol{B}}_Y \boldsymbol{g}^\top\overline{\sigma}_Y^{-1}\widehat{\boldsymbol{B}}_X]^\top$ with $\overline{\sigma}_Y^2\boldsymbol{I}_n = q^{-1}\sum_{k\in K}\mathbb{V}(\widehat{\boldsymbol{\beta}}_{Y_k})$ [2] and similarly for $[\boldsymbol{g}^\top\boldsymbol{B}_Y^* \boldsymbol{g}^\top\boldsymbol{B}_X^*]^\top$. The covariance matrix can be can be partitioned into

$$\boldsymbol{\Sigma}^* = \left[\begin{array}{cc} \boldsymbol{\Sigma}_{YY}^* & \boldsymbol{\Sigma}_{YX}^* \\ \boldsymbol{\Sigma}_{XY}^* & \boldsymbol{\Sigma}_{XX}^* \end{array}\right],$$

where $\boldsymbol{\Sigma}_{XX}^*$, $\boldsymbol{\Sigma}_{YY}^*$ and $\boldsymbol{\Sigma}_{XY}^*$ are the genetic covariances within the exposures, the outcomes and between them, and its inverse ([66], Theorem 8.5.11) into

$$\boldsymbol{\Omega}^* = \boldsymbol{\Sigma}^{*-1} = \left[\begin{array}{cc} \boldsymbol{\Omega}_{YY}^* & -\boldsymbol{\Omega}_{YY}^*\boldsymbol{\Sigma}_{YX}^*\boldsymbol{\Sigma}_{XX}^{*-1} \\ -\boldsymbol{\Sigma}_{XX}^{*-1}\boldsymbol{\Sigma}_{XY}^*\boldsymbol{\Omega}_{YY}^* & \boldsymbol{\Sigma}_{XX}^{*-1} + \boldsymbol{\Sigma}_{XX}^{*-1}\boldsymbol{\Sigma}_{XY}^*\boldsymbol{\Omega}_{YY}^*\boldsymbol{\Sigma}_{YX}^*\boldsymbol{\Sigma}_{XX}^{*-1} \end{array}\right], \tag{8}$$

with $\boldsymbol{\Omega}^* \in \mathcal{P}_\mathcal{D}$, the space of the precision matrices Markov with respect to the DAG $\mathcal{D}$ and $\boldsymbol{\Omega}_{YY}^* = (\boldsymbol{\Sigma}_{YY}^* - \boldsymbol{\Sigma}_{YX}^*\boldsymbol{\Sigma}_{XX}^{*-1}\boldsymbol{\Sigma}_{XY}^*)^{-1}$. However, since by partial ordering $\boldsymbol{\Omega}_{YX}^* = \boldsymbol{\Omega}_{YY}^*\boldsymbol{\Sigma}_{YX}^*\boldsymbol{\Sigma}_{XX}^{*-1} = \boldsymbol{0}$, (8) becomes

$$\boldsymbol{\Omega}^* = \left[\begin{array}{cc} \boldsymbol{\Omega}_{YY}^* & \boldsymbol{0} \\ -\boldsymbol{\Sigma}_{XX}^{*-1}\boldsymbol{\Sigma}_{XY}^*\boldsymbol{\Omega}_{YY}^* & \boldsymbol{\Sigma}_{XX}^{*-1} \end{array}\right].$$

By using $\boldsymbol{\Omega}^*$, Gaussian graphical models [17] can be used to estimate the conditional dependence relationships between the traits in the original (individual-level) data unconfounded by $U$ since genetically predicted outcomes and exposures depend only on the

selected IVs.

Finally, for a given DAG $\mathcal{D}$, the estimand of the causal effect under intervention [67] is

$$\theta_{h,k} = [\boldsymbol{\Sigma}^{*-1}_{\mathrm{fa}(h),\mathrm{fa}(h)} \boldsymbol{\Sigma}^{*}_{\mathrm{fa}(h),k}]_1, \tag{9}$$

where $\boldsymbol{\Sigma}^{*}_{\mathrm{fa}(h),\mathrm{fa}(h)}$ indicates the submatrix of $\boldsymbol{\Sigma}^{*}$ whose rows and columns are $\mathrm{fa}(h)$, $\boldsymbol{\Sigma}^{*}_{\mathrm{fa}(h),k}$ indicates the subvector of $\boldsymbol{\Sigma}^{*}$ whose rows are $\mathrm{fa}(h)$ and the column correspond to the $k$th outcome, and where the subscript indicates the first element of the vector. By using (6) which is a sufficient statistic for $\boldsymbol{\Sigma}^{*}$ after IVW and with $\boldsymbol{V} = \boldsymbol{I}_n$, (9) becomes

$$\theta_{h,k} = [(\boldsymbol{B}^{*\top}_{X_{\mathrm{fa}(h)}} \boldsymbol{B}^{*}_{X_{\mathrm{fa}(h)}})^{-1} \boldsymbol{B}^{*\top}_{X_{\mathrm{fa}(h)}} \boldsymbol{\beta}^{*}_{Y_k}]_1,$$

where $\boldsymbol{B}^{*}_{X_{\mathrm{fa}(h)}}$ are the IVW coefficients of the regressions of each exposure in $\boldsymbol{X}_{\mathrm{fa}(h)}$ on $\boldsymbol{G}$ and $\boldsymbol{\beta}^{*}_{Y_k}$ is the IVW coefficient of a regression of $Y_k$ on $\boldsymbol{G}$. The corresponding estimator coincides with (3).

An important aspect of the MrDAG model is that the genetically predicted values of the outcomes and exposures do not need to be calculated since MrDAG uses as input (6) the sufficient statistic for $\boldsymbol{\Sigma}^{*}$. The only information that would be required from the original (individual-level) data is the LD matrix $\boldsymbol{V}$. However, this information is not necessary in MrDAG summary-level MR design since independent genetic variants are considered after pruning or clumping and thus $\boldsymbol{V} = \boldsymbol{I}_n$.

# MrDAG algorithm

## Markov Equivalent Class, Completed Partially DAGs, Essential Graphs and Partially DAGs

The estimation of a DAG from observational data suffers the known problem of identifiability, *i.e.*, it is not possible to estimate uniquely the underlying true DAG since its conditional independencies can be encoded in several alternative DAGs. This set of DAGs that hold the same conditional independencies is known as Markov Equivalent Class and the best that can be done from observational data is to estimate this class. All DAGs with the same conditional independencies can be represented by a Completed Partially DAG (CPDAG) [68] or Essential Graph (EG) [69]. EGs are Chain Graphs (CGs) whose chain components are decomposable undirected graphs [17]. A CPDAG or EG is a partially directed graph that might contain both directed and undirected edges without directed cycles. Finally, Partially DAG (PDAG) contain both directed and undirected edges and directed cycles might be present.

## Posterior probability of edge inclusion

Technical details of the algorithm for graphical models exploration that we used to develop MrDAG algorithm are presented in [70]. Briefly, it is based on a Markov chain Monte Carlo (MCMC) algorithm devised to explore the space of EGs whose enumeration is infeasible since their number grows super-exponentially with the number of nodes. The EG $\mathcal{G}$ is sampled from a proposal distribution which is accepted with a probability given by a Metropolis-Hastings (M-H) ratio defined to guarantee the convergence of the algorithm to the correct posterior distribution. The key ingredient in the M-H ratio is a closed-form expression for the marginal likelihood $m_{\mathcal{G}}(\text{data})$. This is based on a non-informative prior coupled with a fractional Bayes factor methodology and compatible priors building procedure. In practice, a specific DAG $\mathcal{D}(\mathcal{G})$, which belongs to the Markov Equivalent Class whose unique representative chain graph is the EG $\mathcal{G}$, is proposed and, if accepted, its information stored as an adjacent matrix at each sweep of the MCMC algorithm.

In MrDAG algorithm, we added an acceptance/rejection step to guarantee that $\mathcal{D}(\mathcal{G})$ satisfies the partial ordering that corresponds to the orientation of the edges from the exposures to the outcomes, see Figure 1E. To check the efficiency of this step, we also monitor its acceptance rate. We also included a tempering scheme [71] by considering an annealing parameter $T$ in the M-H ratio to facilitate the convergence of the MCMC algorithm to the target distribution and the exploration of regions of high posterior mass. The temperature $1/T$ exponentiates the M-H ratio and its value increases linearly during the burn-in until $T = 1$ at the end of the burn-in.

Sparsity is enforced by assigning a prior to $\mathcal{G}$ and specifically on $\mathcal{G}^U$, the skeleton of $\mathcal{G}$ which contains the same edges of $\mathcal{G}$ but without orientation

$$\mathcal{G}^U_{(l)} \mid \pi^{\text{edge}} \overset{\text{i.i.d.}}{\sim} \text{Ber}(\pi^{\text{edge}}), \quad l = 1, \ldots, (q+p)(q+p-1)/2, \tag{10}$$

where $\mathcal{G}^U_{(l)}$ is the $l$th element of the vectorized lower triangular part of the adjacency matrix of $\mathcal{G}^U$ and $(q+p)(q+p-1)/2$ is the maximum number of edges in an EG on $q+p$ nodes.

The posterior distribution of $\mathcal{G}$ is

$$\mathbb{P}(\mathcal{G} \mid \text{data}) = \frac{m_{\mathcal{G}}(\text{data})\,\mathbb{P}(\mathcal{G})}{\sum_{\mathcal{G} \in \mathcal{S}_{q+p}} m_{\mathcal{G}}(\text{data})\,\mathbb{P}(\mathcal{G})} \tag{11}$$

with $\mathcal{S}_{q+p}$ the set of all EGs with $q+p$ nodes. The posterior probability of edge inclusion (PPeIs) is defined as

$$\begin{aligned}\mathbb{P}_{z \to v}(\text{data}) &= \textstyle\sum_{\mathcal{G} \in \mathcal{S}_{z \to v}} \mathbb{P}(\mathcal{G} \mid \text{data}) \\ &\approx \tfrac{1}{S} \textstyle\sum_{s=1}^{S} \mathbb{I}_{\mathcal{D}(\mathcal{G}^{(s)})}(\mathcal{D}(\mathcal{G}^{(s)}_{z \to v})),\end{aligned} \tag{12}$$

where $\mathcal{S}_{z \to v}$ is the set of EGs containing the directed edge $z \to v$, $S$ is the number of

sweeps after burn-in and $\mathbb{I}_{\mathcal{D}(\mathcal{G}^{(s)})}(\mathcal{D}(\mathcal{G}^{(s)}_{z\to v}))$ is the indicator function that is equal to one if the specific DAG considered at the $s$th sweep $\mathcal{D}(\mathcal{G}^{(s)})$ contains the directed edge $z \to v$ and zero otherwise.

Note that, although MrDAG explores the space of EGs and stores a specific DAG that belongs to the sampled EG, the graphs obtained by thresholding the PPeIs might give rise to a PDAG [70].

## Bayesian causal effects estimation

Here, we summarise the results reported in [67] that we employed to derive the Bayesian estimation of the causal effects under unmeasured confounders.

Let's rewrite (7) as

$$[\boldsymbol{g}^\top \widehat{\boldsymbol{B}}^*_Y \boldsymbol{g}^\top \widehat{\boldsymbol{B}}^*_X]^\top \mid \boldsymbol{\Sigma}^*_{\mathcal{D}} \sim \mathrm{N}_{q+p}([\boldsymbol{g}^\top \boldsymbol{B}^*_Y \boldsymbol{g}^\top \boldsymbol{B}^*_X]^\top, \boldsymbol{\Sigma}^*_{\mathcal{D}}),$$

where $\boldsymbol{\Sigma}^*_{\mathcal{D}} \in \mathcal{C}_{\mathcal{D}}$, the space of s.p.d. covariance matrices Markov with respect to $\mathcal{D}$. In the following, for ease of notation, we refer to $[\boldsymbol{g}^\top \widehat{\boldsymbol{B}}^*_Y \boldsymbol{g}^\top \widehat{\boldsymbol{B}}^*_X]^\top$ as the "data" and we also drop the subscript $\mathcal{D}$.

Let $\boldsymbol{\Omega}^* = \boldsymbol{\Sigma}^{*^{-1}} = \boldsymbol{L}^* \boldsymbol{D}^{*^{-1}} \boldsymbol{L}^{*\top}$ be the modified Cholesky decomposition of the precision $\boldsymbol{\Omega}^*$. The DAG Cholesky parametrization of $\boldsymbol{\Omega}^*$ is given by the node-parameters $\boldsymbol{\omega}^*_l = (D^*_{ll}, \boldsymbol{L}^*_{\mathrm{fa}(l)})$, $l = 1, \ldots, q + p$, with

$$D^*_{ll} = \boldsymbol{\Sigma}^*_{ll\mid\mathrm{pa}(l)}, \quad \boldsymbol{L}^*_{\mathrm{fa}(l)} = -\boldsymbol{\Sigma}^{*^{-1}}_{\mathrm{pa}(l)} \boldsymbol{\Sigma}^*_{\mathrm{fa}(l)},$$

where $\boldsymbol{\Sigma}^*_{\mathrm{pa}(l)}$ indicates the submatrix of $\boldsymbol{\Sigma}^*$ whose rows and columns are $\mathrm{pa}(l)$.

For a given DAG $\mathcal{D}$, [67] derive the posterior distribution of $\boldsymbol{\omega}^*_l$, $l = 1, \ldots, q + p$, in an objective Bayes framework which has the advantage of not depending on priors hyperparameters. In turn, the posterior draws of the Cholesky parameters $\boldsymbol{\omega}^*_l$ provide posterior draws from $(\boldsymbol{\Omega}^* \mid \text{data}) = (\boldsymbol{L}^* \boldsymbol{D}^{*^{-1}} \boldsymbol{L}^{*\top} \mid \text{data})$ and finally, by using (9), posterior samples of the causal effects between the exposures and the outcomes.

In contrast to frequentist approaches [72] where, for an estimated EG $\mathcal{G}$, the causal effects are calculated averaging over all (if numerical feasible) or a subset of DAGs within the Markov Equivalent Class $\mathcal{G}$, here, we also consider the uncertainty related to the estimation of the EGs. Let $\{\mathcal{G}_v, v = 1, \ldots, V\}$ the set of unique visited EGs by MrDAG. Based on (11), the posterior probability of $\mathcal{G}_v$ can be approximated by

$$\mathbb{P}(\mathcal{G}_v \mid \text{data}) \approx \frac{m_{\mathcal{G}_v}(\text{data})\,\mathbb{P}(\mathcal{G}_v)}{\sum_{v=1}^V m_{\mathcal{G}_v}(\text{data})\,\mathbb{P}(\mathcal{G}_v)}. \tag{13}$$

Averaging over the unique visited EGs, the posterior causal effect under intervention in

the $h$th exposures on the $k$th outcomes is

$$
\begin{aligned}
\theta_{h,k} \mid \text{data} &= \sum_{v=1}^{V} \mathbb{E}(\theta_{h,k}(\mathcal{D}(\mathcal{G}_v)) \mid \text{data}, \mathcal{G}_v) \, \mathbb{P}(\mathcal{G}_v \mid \text{data}), \quad h \in J, k \in K, \\
&\approx \frac{1}{S} \sum_{s=1}^{S} \theta_{h,k}(\mathcal{D}(\mathcal{G}^{(s)})),
\end{aligned}
\tag{14}
$$

where $\mathcal{G}_v$ is one of the unique EGs visited during the MCMC, $\mathbb{E}(\theta_{h,k}(\mathcal{D}(\mathcal{G}_v)) \mid \text{data}, \mathcal{G}_v)$ is the posterior expectation of the causal effect given $\mathcal{G}_v$, *i.e.*, over *all* $\mathcal{D}(\mathcal{G}_v)$, $\mathbb{P}(\mathcal{G}_v \mid \text{data})$ is defined in (13) and $\theta_{h,k}(\mathcal{D}(\mathcal{G}^{(s)})$ is the posterior causal effect conditioned on the recorded DAG at the $s$th sweep.

Finally, by a suitable modification of (14), credible intervals of the causal effects between the exposures and outcomes can be derived.

## Simulation study

We share several aspects of the simulation study with [2]. It is formulated in a two-sample summary-level MR design, where $N = 100,000$ independent individuals are simulated, of which $N_Y = 50,000$ are used to compute the genetic associations with the exposures and $N_X = 50,000$ to compute the genetic associations with the outcomes. Thus, we assume that the quantitative exposures $X_j$, $j \in J = \{1, \ldots, p\}$, and the quantitative responses $Y_k$, $k \in K = \{1, \ldots, q\}$, are measured on the same individuals $N_X$ and $N_Y$, respectively, with 100% sample overlap, but independent of each other.

In all simulated scenarios, we consider $p = 15$ exposures, $q = 5$ outcomes and $n = 100$ independent genetic variants as IVs. Genotypes for the $i$th genetic variant and each individual $\ell$ are simulated independently according to a binomial distribution with minor allele frequency (MAF) equal to 0.05, *i.e.*, $g_{\ell,i} \overset{\text{i.i.d.}}{\sim} \text{Bin}(2, 0.05)$, $\ell \in L = \{1, \ldots, N\}$, $i \in I = \{1, \ldots, n\}$. The resulting matrix of genotypes $\boldsymbol{G}$ is split into two equally sized groups, $\boldsymbol{G}_X$ and $\boldsymbol{G}_Y$, of dimension $N_X \times n$ and $N_Y \times n$, respectively. Thus, no IVW is needed in the simulation study given that the same MAF at 5% is used to simulate the genotypes.

Overall, the data generation process consists of two stages. In the first stage, the raw data for the exposures $\boldsymbol{X}$ and the outcomes $\boldsymbol{Y}$ are simulated. Then, in the second stage, summary-level statistics are obtained as the linear regression coefficients $\widehat{\beta}_{i,j}^{X}$ from a univariable linear regression in which the $j$th exposure is regressed on the $i$th genetic variant in sample one and the linear regression coefficients $\widehat{\beta}_{i,k}^{Y}$ from a univariable linear regression in which the $k$th outcome is regressed on the $i$th genetic variant in sample two.

In the following, we detail each stage and how we simulate the quantities involved. We start with the first stage which is divided into two steps.

- In the first step, the exposures are generated as follows

$$
\boldsymbol{x}_j = \boldsymbol{G}_X \boldsymbol{\beta}_{X_j} + \psi_X \boldsymbol{u}_X + \boldsymbol{\epsilon}_{X_j}, \quad j \in J,
\tag{15}
$$

where $\boldsymbol{G}_X$ and $\boldsymbol{u}_X$ are the genotypes of the $n$ IVs and the values of the confounder $U$ measured on the same $N_X$ individuals, respectively, and where $\boldsymbol{\beta}_{X_j}$ and $\psi_X$ are the corresponding genetic and confounding effects. $\boldsymbol{\epsilon}_{X_j} \sim \mathrm{N}_{N_X}(\boldsymbol{0}, h_{X_j}\boldsymbol{I}_{N_X})$ with $h_{X_j}$ the $j$th diagonal element of the $(p \times p)$-dimensional matrix

$$\boldsymbol{H}_X = \frac{1 - v_X}{v_X}(\boldsymbol{G}_X\boldsymbol{B}_X + \psi_X\boldsymbol{u}_X\boldsymbol{1}_p^\top)^\top(\boldsymbol{G}_X\boldsymbol{B}_X + \psi_X\boldsymbol{u}_X\boldsymbol{1}_p^\top), \qquad (16)$$

where $v_X$ is the desired level of heritability, or how much variation $\boldsymbol{G}$ can explain of $X_j$, fixed at 10% for all exposures and in all simulated scenarios. In (16), $\boldsymbol{B}_X = \{\boldsymbol{\beta}_{X_j}\}_{j \in J}$ is an $(n \times p)$-dimensional matrix of the effects of the genetic variants on the exposures.

The confounder $U$ is drawn from a multivariate standard Gaussian distribution, i.e., $\boldsymbol{u} \sim \mathrm{N}_N(\boldsymbol{0}, \boldsymbol{I}_N)$ and, then, split into two equally sized vectors $\boldsymbol{u}_X$ and $\boldsymbol{u}_Y$ with effect $\psi_X$ impacting all exposures and $\psi_Y$ effecting all outcomes.

The effects $\boldsymbol{\beta}_{X_j}$ of the $n$ genetic variants on the $j$th exposure are drawn following [70]. We randomly generate a topologically ordered DAG among the $p$ exposures with a probability of edge inclusion $p_X^{\mathrm{edge}} = 2/(p-1)$ using the function `randomDAG()` in the R package *pcalg* [27]. Thus, the resulting DAG implies the following system of equations [18]

$$\boldsymbol{\beta}_{X_j} = \sum_{h \in \mathrm{pa}(j)} \gamma_{h,j}^X \boldsymbol{\beta}_{X_h} + \boldsymbol{\epsilon}_{X_j} \qquad (17)$$

with $\boldsymbol{\epsilon}_{X_j} \sim \mathrm{N}(\boldsymbol{0}, \boldsymbol{I}_n)$. For each $j \in J$, the effect within the exposures $\gamma_{h,j}^X$ are uniformly chosen in the interval $[-1.1r_X, -0.9r_X] \cup [0.9r_X, 1.1r_X]$. This construction procedure for $\boldsymbol{\beta}_{X_j}$ corresponds to the simulated scenario that we call "DAG$_X$", i.e., Directed Acyclic Graph within $\boldsymbol{X}$, which, in turn, is paired with two different simulated scenarios for the effects $\boldsymbol{\beta}_{Y_k}$ described in the second step (first stage) of the simulation study.

We also simulate the effects $\boldsymbol{\beta}_{X_j}$ following [2]. Specifically, we simulate $\boldsymbol{\beta}_{X_j} \sim \mathrm{N}_n(\boldsymbol{0}, \boldsymbol{R}_X)$, where $\boldsymbol{R}_X$ is the $(p \times p)$-dimensional Toeplitz matrix with $r_X^{|j-j'|}$ for $j, j' \in J$. The matrix $\boldsymbol{R}_X$ implies a tridiagonal sparse inverse correlation matrix $\boldsymbol{\Omega}_X = \boldsymbol{R}_X^{-1}$. The interpretation of non-zero elements of $\boldsymbol{\Omega}_X$ coincides with the effects simulated in (17). We call this second scenario for the effects of the genotypes on the exposures "UG$_X$", i.e., Undirected Graph within $\boldsymbol{X}$.

In both simulated scenarios for $\boldsymbol{X}$, we use different levels of $r_X$, ranging from independence to a strong dependence, i.e., $r_X = \{0, 0.2, 0.4, 0.6, 0.8\}$, where $r_X = 0.6$ represents a medium dependence between the genetic associations with the exposures. We use this value in the figures presented in Section 'Simulation study'.

- In the second step (first stage) of the simulation study, the outcomes are generated

on another independent set of $N_Y$ individuals based on the following set of equations

$$\boldsymbol{y}_k = \boldsymbol{X}\boldsymbol{\theta}_k + \sum_{h\in\mathrm{pa}(k)} \gamma^Y_{h,k}\boldsymbol{y}_h + \psi_Y\boldsymbol{u} + \boldsymbol{\epsilon}_{Y_k}, \quad k \in K, \tag{18}$$

where $\boldsymbol{X}$ is the $(N_X \times p)$-dimensional matrix of exposures simulated using (15), $\boldsymbol{\theta}_k = (\theta_{1k}, \ldots, \theta_{pk})^\top$ is $p$-dimensional (sparse) vector the causal effects from the exposures to the $k$th outcome and where $\psi_Y$ is the effect of the confounder $U$ on the outcomes. $\boldsymbol{\epsilon}_{Y_k} \sim \mathrm{N}_{N_Y}(\boldsymbol{0}, h_{Y_k}\boldsymbol{I}_{N_Y})$ with $h_{Y_k}$ the $k$th diagonal element of the $(q \times q)$-dimensional matrix

$$\boldsymbol{H}_Y = \frac{1-v_Y}{v_Y}(\boldsymbol{X}\boldsymbol{\theta}_k + \sum_{h\in\mathrm{pa}(k)} \gamma^Y_{h,k}\boldsymbol{y}_h + \psi_Y\boldsymbol{u} + \boldsymbol{\epsilon}_{Y_k})^\top$$
$$(\boldsymbol{X}\boldsymbol{\theta}_k + \sum_{h\in\mathrm{pa}(k)} \gamma^Y_{h,k}\boldsymbol{y}_h + \psi_Y\boldsymbol{u} + \boldsymbol{\epsilon}_{Y_k}),$$

where $v_Y$ is the desired level of the proportion of variance explained, fixed at 25% for all outcomes and in all simulated scenarios.

In (18), the term $\sum_{h\in\mathrm{pa}(k)} \gamma^Y_{h,k}\boldsymbol{Y}_h$ depends on a randomly generated topologically ordered DAG among the $q$ outcomes with probability of edge inclusion $p_Y^{\mathrm{edge}} = 1/(q-1)$. For each $k \in K$, the effects within the outcomes $\gamma^Y_{h,k}$ are uniformly drawn in the interval $[0.9m_Y, 1.1m_Y]$. In analogy with the first step, we call this scenario "DAG$_Y$", i.e., Directed Acyclic Graph within $\boldsymbol{Y}$.

We also simulate a simplified scenario where

$$\boldsymbol{y}_k = \gamma^Y_{h,k}\boldsymbol{y}_h + \psi_Y\boldsymbol{u} + \boldsymbol{\epsilon}_{Y_k}, \tag{19}$$

i.e., a randomly selected outcome $k$ is completed mediated by another randomly selected response chosen between the remaining ones. We call this scenario "Med$_Y$", i.e., complete mediation of an outcome, since in the previous scenario "DAG$_Y$" partial mediations [73] are likely simulated, while here we exclude this case. In this second simulated scenario for the outcomes, the matrix $\boldsymbol{H}_Y$ is calculated according to (19). Moreover, we use different levels of $m_Y$, ranging from small to a strong level of (partial or complete) mediation, i.e., $m_Y = \{0.25, 0.50, 0.75, 1, 1.5, 2\}$, where $m_Y = 1$ represents a medium (partial or complete) mediation effect. We use this value in the figures presented in Section 'Simulation study'.

Finally, the causal effects $\boldsymbol{\theta}_k$ are drawn independently from a multivariate Gaussian distribution, i.e., $\boldsymbol{\theta}_k \sim \mathrm{N}_p(\boldsymbol{0}, \boldsymbol{I}_p)$.

In both simulated scenarios for $\boldsymbol{Y}$, we consider a $(q \times p)$-dimensional sparse matrix of causal effects $\boldsymbol{\Theta} = \{\boldsymbol{\theta}_k\}_{k\in K}$, where 30 cells of the matrix are non-zero and where several exposures are either *shared* or *distinct* for the outcomes. Specifically, we select at random the same proportion of cells in the matrix $\boldsymbol{\Theta}$ and assign them the

simulated values, while the other cells are set to zero.

After the first stage, four scenarios are created by combining the simulations for $\boldsymbol{X}$ and $\boldsymbol{Y}$: (i) "UndG$_X$-Med$_Y$", *i.e.*, undirected graph within $\boldsymbol{X}$ and complete mediation of an outcome in $\boldsymbol{Y}$; (ii) DAG$_X$-Med$_Y$, *i.e.*, topologically ordered DAG within $\boldsymbol{X}$ and complete mediation of a response within $\boldsymbol{Y}$; (iii) UndG$_X$-DAG$_Y$, *i.e.*, undirected graph within $\boldsymbol{X}$ and topologically ordered DAG within $\boldsymbol{Y}$; (iv) DAG$_X$-DAG$_Y$, topologically ordered DAGs within $\boldsymbol{X}$ and $\boldsymbol{Y}$. In (ii) and (iv) the overall DAGs, obtained by combining different simulation patterns for $\boldsymbol{X}$ and $\boldsymbol{Y}$, are fully oriented while in (i) and (iii) they are partially oriented.

After creating the data at the individual level, in the second stage, we compute the summary-level statistics from the two independent groups of individuals. The input data for the simulation study are the summary-level statistics $\widehat{\boldsymbol{B}}_X = \{\widehat{\beta}_{i,j}^X\}_{i\in I, j\in J}$, an $(n \times p)$-dimensional matrix, and $\widehat{\boldsymbol{B}}_Y = \{\widehat{\beta}_{i,k}^Y\}_{i\in I, k\in K}$, an $(n \times q)$-dimensional matrix, derived from a univariable linear regression model, where each genetic variant $G_i$ is regressed against each exposure $X_j$ and each outcome $Y_k$ at-a-time.

## Real data application: Pre-processing and data preparation

The first step of the data processing merges the summary-level data (beta regression coefficients, their standard errors and associated $p$-values) of all exposures by their unique "rs" identifier and aligns the effect direction of the genetic associations with each exposure according to the same effect allele. As IVs, we select the genetic variants which are associated with any of the exposures at genome-wide significance (minimum $p$-value $< 5 \times 10^{-8}$ across all exposures). Next, we merge the genetic variants selected as IVs with the outcome data by their unique "rs" identifier and align the effect direction of the genetic associations with each outcome according to the same effect allele. Finally, we clump the genetic variants to be independent at $r^2 < 0.01$ using a European reference panel [74]. This results in $n = 708$ independent genetic variants selected as IVs. See Supplementary Table 1 for the description of the summary-level statistics, the data sources, the number of non-unique IVs which were genome-wide significant for each exposure along with the contribution (%) of each exposure on the selected IVs.

Finally, we perform reverse causation using the same traits with mental health phenotypes as exposures and lifestyle and behavioural traits as outcomes. We apply the same procedure described above resulting in 470 IVs. See Supplementary Table 1 for details regarding the number of non-unique IVs which were genome-wide significant for each exposure along with the contribution (%) of each exposure on the selected IVs.

# Data availability

Data sources are presented in Supplementary Information with associated URL links. Social Science Genetic Association Consortium (SSGAC) summary-level statistics are available through a standard registration procedure (`https://thessgac.com/register/`).

# Code availability

The Mendelian randomization with Directed Acyclic Graph learning R package `MrDAG` is freely available on `https://github.com/lb664/MrDAG/`. It includes the data of the real data examples and how to run the algorithm. Post-processing routines to estimate the posterior causal effects presented in the manuscript are also included along with the Posterior Probability of Edge Inclusion.

# Declaration of interests

The authors do not have competing interests.

# Acknowledgements

# Author contributions

Conceptualization: L.B., V.Z., D.G.; Methodology: L.B, V.Z.; Formal Analysis: V.Z., L.B.; Resources: (all collaborators); Data curation: V.Z., D.G.; Writing – original draft: V.Z., L.B., T.C., G.D., N.C.; Visualization: V.Z., L.B.; Supervision: L.B., D.G.

# Funding

MRC, Alzheimer's Society and Alzheimer's Research UK, and by the NIHR Cambridge Biomedical Research Centre (NIHR203312) (L.B.). The views expressed are those of the authors and not necessarily those of the NIHR or the Department of Health and Social Care.

# Competing interests

The authors declare no competing interests.

# Additional information

Supplementary Information includes technical details regarding the identifiability and the estimation of causal effects for a given DAG, the consistency of the effects of the regressions of the exposures and the outcomes on the selected genetic variants along with figures and tables to support the results of the simulation study and the real data analysis.

# References

[1] Zuber, V., Colijn, J. M., Klaver, C. & Burgess, S. Selecting likely causal risk factors from high-throughput experiments using multivariable Mendelian randomization. *Nature Communications* **11**, 29 (2020). URL https://doi.org/10.1038/s41467-019-13870-3.

[2] Zuber, V. *et al.* Multi-response Mendelian randomization: Identification of shared and distinct exposures for multimorbidity and multiple related disease outcomes. *American Journal of Human Genetics* **110**, 1177–1199 (2023). URL https://doi.org/10.1016/j.ajhg.2023.06.005.

[3] Sanderson, E., Davey Smith, G., Windmeijer, F. & Bowden, J. An examination of multivariable Mendelian randomization in the single-sample and two-sample summary data settings. *International Journal of Epidemiology* **48**, 713–727 (2019). URL https://doi.org/10.1093/ije/dyy262.

[4] Hernán, M. A. & Robins, J. M. Instruments for causal inference. An epidemiologist's dream? *Epidemiology* **17**, 360–372 (2006). URL https://doi.org/10.1097/01.ede.0000222409.00878.37.

[5] Didelez, V. & Sheehan, N. Mendelian randomization as an instrumental variable approach to causal inference. *Statistical Methods in Medical Research* **16**, 309–330 (2007). URL https://doi.org/10.1177/0962280206077743.

[6] Drton, M. & Maathuis, M. H. Structure learning in graphical modeling. *Annual Review of Statistics and Its Application* **4**, 365–393 (2017). URL https://doi.org/10.1146/annurev-statistics-060116-053803.

[7] Pearl, J. *Causality: Models, Reasoning and Inference* (Cambridge University Press, Cambridge, 2009), 2nd edn.

[8] GBD 2019 Mental Disorders Collaborators. Global, regional, and national burden of 12 mental disorders in 204 countries and territories, 1990–2019: A systematic analysis for the Global Burden of Disease Study 2019. *Lancet Psychiatry* **9**, 137–150 (2022). URL https://doi.org/10.1016/s2215-0366(21)00395-3.

[9] Arias, D., Saxena, S. & Verguet, S. Quantifying the global burden of mental disorders and their economic value. *eClinicalMedicine* **54**, 101675 (2022). URL https://doi.org/10.1016/j.eclinm.2022.101675.

[10] Tyrer, P. A comparison of DSM and ICD classifications of mental disorder. *Advances in Psychiatric Treatment* **20**, 280–285 (2014). URL https://doi.org/10.1192/apt.bp.113.011296.

[11] Saxe, G. N., Bickman, L., Ma, S. & Aliferis, C. Mental health progress requires causal diagnostic nosology and scalable causal discovery. *Frontiers in Psychiatry* **13**, 898789 (2022). URL https://doi.org/10.3389/fpsyt.2022.898789.

[12] Arango, C. *et al.* Risk and protective factors for mental disorders beyond genetics: An evidence-based atlas. *World Psychiatry* **20**, 417–436 (2021). URL https://doi.org/10.1002/wps.20894.

[13] Leichsenring, F., Steinert, C., Rabung, S. & Ioannidis, J. P. A. The efficacy of psychotherapies and pharmacotherapies for mental disorders in adults: An umbrella review and meta-analytic evaluation of recent meta-analyses. *World Psychiatry* **21**, 133–145 (2022). URL https://doi.org/10.1002/wps.20941.

[14] Saccaro, L. F., Gasparini, S. & Rutigliano, G. Applications of Mendelian randomization in psychiatry: A comprehensive systematic review. *Psychiatric Genetics* **32**, 199–213 (2022). URL https://doi.org/10.1097/ypg.0000000000000327.

[15] Smith, G. D. & Ebrahim, S. 'Mendelian randomization': Can genetic epidemiology contribute to understanding environmental determinants of disease? *International Journal of Epidemiology* **32**, 1–22 (2003). URL https://doi.org/10.1093/ije/dyg070.

[16] Thomas, D. C. & Conti, D. V. Commentary: The concept of 'Mendelian Randomization'. *International Journal of Epidemiology* **33**, 21–25 (2004). URL https://doi.org/10.1093/ije/dyh048.

[17] Lauritzen, S. L. *Graphical Models* (Clarendon Press, Oxford, New York, 2004). Repr. with corrections.

[18] Peters, J. & Bühlmann, P. Identifiability of Gaussian structural equation models with equal error variances. *Biometrika* **101**, 219–228 (2013). URL https://doi.org/10.1093/biomet/ast043.

[19] Perković, E., Kalisch, M. & Maathuis, M. H. Interpreting and using CPDAGs with background knowledge. In *Proceedings UAI* (AUAI Press, 2017). URL http://auai.org/uai2017/proceedings/papers/120.pdf.

[20] Didelez, V., Meng, S. & Sheehan, N. A. Assumptions of IV methods for observational epidemiology. *Statistical Science* **25**, 22–40 (2010). URL https://doi.org/10.1214/09-sts316.

[21] Frot, B., Nandy, P. & Maathuis, M. H. Robust causal structure learning with some hidden variables. *Journal of the Royal Statistical Society: Series B* **81**, 459–487 (2019). URL https://doi.org/10.1111/rssb.12315.

[22] Grotzinger, A. D. *et al.* Genomic structural equation modelling provides insights into the multivariate genetic architecture of complex traits. *Nature Human Behaviour* **3**, 513–525 (2019). URL https://doi.org/10.1038/s41562-019-0566-x.

[23] Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nature Genetics* **47**, 1236–1241 (2015). URL https://doi.org/10.1038/ng.3406.

[24] Shi, H., Mancuso, N., Spendlove, S. & Pasaniuc, B. Local genetic correlation gives insights into the shared genetic architecture of complex traits. *American Journal of Human Genetics* **101**, 737–751 (2017). URL https://doi.org/10.1016/j.ajhg.2017.09.022.

[25] Burgess, S. & Bowden, J. Integrating summarized data from multiple genetic variants in Mendelian randomization: Bias and coverage properties of inverse-variance weighted methods (2015). URL https://doi.org/10.48550/ARXIV.1512.04486. Unpublished manuscript.

[26] Wu, Y., Kang, H. & Ye, T. Debiased multivariable Mendelian randomization (2024). URL https://doi.org/10.48550/ARXIV.2402.00307. Preprint, 2402.00307.

[27] Kalisch, M., Mächler, M., Colombo, D., Maathuis, M. H. & Bühlmann, P. Causal inference using graphical models with the R package pcalg. *Journal of Statistical Software* **47** (2012). URL https://doi.org/10.18637/jss.v047.i11.

[28] Badsha, M. B., Martin, E. A. & Fu, A. Q. MRPC: An R package for inference of causal graphs. *Frontiers in Genetics* **12** (2021). URL https://doi.org/10.3389/fgene.2021.651812.

[29] Spirtes, P., Glymour, C. & Scheines, R. *Causation, Prediction and Search* (MIT Press, Cambridge, MA, 2000), 2nd edn.

[30] Rahman, S., Khare, K., Michailidis, G., Martínez, C. & Carulla, J. Estimation of Gaussian Directed Acyclic Graphs using partial ordering information with applications to DREAM3 networks and dairy cattle data. *The Annals of Applied Statistics* **17**, 929–960 (2023). URL https://doi.org/10.1214/22-aoas1636.

[31] Raghupathi, V. & Raghupathi, W. The influence of education on health: An empirical assessment of OECD countries for the period 1995–2015. *Archives of Public Health* **78** (2020). URL https://doi.org/10.1186/s13690-020-00402-5.

[32] Amin, V., Fletcher, J. M., Lu, Q. & Song, J. Re-examining the relationship between education and adult mental health in the UK: A research note. *Economics of Education Review* **93**, 102354 (2023). URL https://doi.org/10.1016/j.econedurev.2023.102354.

[33] Davies, N. M., Dickson, M., Smith, G. D., van den Berg, G. J. & Windmeijer, F. The causal effects of education on health outcomes in the UK Biobank. *Nature Human Behaviour* **2**, 117–125 (2018). URL `https://doi.org/10.1038/s41562-017-0279-y`.

[34] Dardani, C. *et al.* Is genetic liability to ADHD and ASD causally linked to educational attainment? *International Journal of Epidemiology* **50**, 2011–2023 (2022). URL `https://doi.org/10.1093/ije/dyab107`.

[35] Verhoef, E. *et al.* Discordant associations of educational attainment with ASD and ADHD implicate a polygenic form of pleiotropy. *Nature Communications* **12**, 6534 (2021). URL `https://doi.org/10.1038/s41467-021-26755-1`.

[36] Cai, J. *et al.* Socioeconomic status, individual behaviors and risk for mental disorders: A Mendelian randomization study. *European Psychiatry* **65**, e28 (2022). `https://doi.org/10.1192/j.eurpsy.2022.18`.

[37] Lloyd, E. C., Reed, Z. E. & Wootton, R. E. The absence of association between anorexia nervosa and smoking: Converging evidence across two studies. *European Child and Adolescent Psychiatry* **32**, 1229–1240 (2023). URL `https://doi.org/10.1007/s00787-021-01918-z`.

[38] Kari, J. T. *et al.* Education leads to a more physically active lifestyle: Evidence based on Mendelian randomization. *Scandinavian Journal of Medicine & Science in Sports* **30**, 1194–1204 (2020). URL `/https://doi.org/10.1111/sms.13653`.

[39] Gage, S. H., Bowden, J., Smith, G. D. & Munafo, M. R. Investigating causality in associations between education and smoking: A two-sample Mendelian randomization study. *International Journal of Epidemiology* **47**, 1131–1140 (2018). URL `https://doi.org/10.1101/184218`.

[40] Wootton, R. E. *et al.* Evidence for causal effects of lifetime smoking on risk for depression and schizophrenia: A Mendelian randomisation study. *Psychological Medicine* **50**, 2435–2443 (2020). URL `https://doi.org/10.1017/S0033291719002678`.

[41] Treur, J. L. *et al.* Investigating causality between liability to ADHD and substance use, and liability to substance use and ADHD risk, using Mendelian randomization. *Addiction Biology* **26**, e12849 (2021). URL `https://doi.org/10.1111/adb.12849`.

[42] West, A. B. *et al.* A systematic review of physical activity, sedentary behavior, and substance use in adolescents and emerging adults. *Translational Behavioral Medicine* **10**, 1155–1167 (2020). URL `https://doi.org/10.1093/tbm/ibaa008`.

[43] Iob, E. *et al.* Testing the causal relationships of physical activity and sedentary behaviour with mental health and substance use disorders: a Mendelian randomisation study. *Molecular Psychiatry* **28**, 3429–3443 (2023). URL `https://doi.org/10.1038/s41380-023-02133-9`.

[44] Reed, Z. E., Wootton, R. E. & Munafo, M. R. Using Mendelian randomisation to explore the gateway hypothesis: Possible causal effects of smoking initiation and alcohol consumption on substance use outcomes. *Addiction* **117**, 741–750 (2022). URL `https://doi.org/10.1111/add.15673`.

[45] Lohr, J. B. & Flynn, K. Smoking and schizophrenia. *Schizophrenia Research* **8**, 93–102 (1992). URL https://doi.org/10.1016/0920-9964(92)90024-Y.

[46] Kendler, K. S., Lönn, S. L., Sundquist, J. & Sundquist, K. Smoking and schizophrenia in population cohorts of Swedish women and men: A prospective co-relative control study. *American Journal of Psychiatry* **172**, 1092–1100 (2015). URL https://doi.org/10.1176/appi.ajp.2015.15010126.

[47] Kendler, K. S., Ohlsson, H., Sundquist, J. & Sundquist, K. Risk for mood, anxiety, and psychotic disorders in individuals at high and low genetic liability for bipolar disorder and major depression. *JAMA Psychiatry* **79**, 1102 (2022). URL https://doi.org/10.1001/jamapsychiatry.2022.2873.

[48] Richards, A. L. *et al.* Genetic liabilities differentiating bipolar disorder, schizophrenia, and major depressive disorder, and phenotypic heterogeneity in bipolar disorder. *JAMA Psychiatry* **79**, 1032 (2022). URL https://doi.org10.1371/j10.1001/jamapsychiatry.2022.2594.

[49] Peyre, H. *et al.* Combining multivariate genomic approaches to elucidate the comorbidity between autism spectrum disorder and attention deficit hyperactivity disorder. *Journal of Child Psychology and Psychiatry* **62**, 1285–1296 (2021). URL https://doi.org/10.1111/jcpp.13479.

[50] Vermeulen, J. M. *et al.* Smoking and the risk for bipolar disorder: Evidence from a bidirectional Mendelian randomisation study. *The British Journal of Psychiatry* **218**, 88–94 (2019). URL https://doi.org/10.1192/bjp.2019.202.

[51] Cai, N. *et al.* Minimal phenotyping yields genome-wide association signals of low specificity for major depression. *Nature Genetics* **52**, 437–447 (2020). URL https://doi.org/10.1038/s41588-020-0594-5.

[52] Cai, N., Choi, K. W. & Fried, E. I. Reviewing the genetics of heterogeneity in depression: Operationalizations, manifestations and etiologies. *Human Molecular Genetics* **29**, R10–R18 (2020). URL https://doi.org/10.1093/hmg/ddaa115.

[53] Forbes, M. K. *et al.* Elemental psychopathology: Distilling constituent symptoms and patterns of repetition in the diagnostic criteria of the DSM-5. *Psychological Medicine* **54**, 886–894 (2023). URL https://doi.org/10.1017/s0033291723002544.

[54] Kendrick, T. *et al.* Management of depression in adults: summary of updated NICE guidance. *British Medical Journal* **378**, o1557 (2022).

[55] Glymour, C., Zhang, K. & Spirtes, P. Review of causal discovery methods based on graphical models. *Frontiers in Genetics* **10** (2019). URL https://doi.org/10.3389/fgene.2019.00524.

[56] Anderson, E. L. *et al.* Education, intelligence and Alzheimer's disease: Evidence from a multivariable two-sample Mendelian randomization study. *International Journal of Epidemiology* **49**, 1163–1172 (2020). URL https://doi.org/10.1093/ije/dyz280.

[57] van Amsterdam, J., van der Velde, B., Schulte, M. & van den Brink, W. Causal factors of increased smoking in ADHD: A systematic review. *Substance Use & Misuse* **53**, 432–445 (2017). URL https://doi.org/10.1080/10826084.2017.1334066.

[58] Green, R., Baker, N. L., Ferguson, P. L., Hashemi, D. & Gray, K. M. ADHD symptoms and smoking outcomes in a randomized controlled trial of varenicline for adolescent and young adult tobacco cessation. *Drug and Alcohol Dependence* **244**, 109798 (2023). URL https://doi.org/10.1016/j.drugalcdep.2023.109798.

[59] Maathuis, M. H., Kalisch, M. & Bühlmann, P. Estimating high-dimensional intervention effects from observational data. *The Annals of Statistics* **37**, 3133–3164 (2009). URL https://doi.org/10.1214/09-aos685.

[60] Badsha, M. B. & Fu, A. Q. Learning causal biological networks with the principle of Mendelian randomization. *Frontiers in Genetics* **10** (2019). URL https://doi.org/10.3389/fgene.2019.00460.

[61] Hoeting, J. A., Madigan, D., Raftery, A. E. & Volinsky, C. T. Bayesian model averaging: A tutorial (with comments by M. Clyde, David Draper and E. I. George, and a rejoinder by the authors). *Statistical Science* **14**, 382–417 (1999). URL https://doi.org/10.1214/ss/1009212519.

[62] Sanderson, E., Davey Smith, G., Bowden, J. & Munafò, M. R. Mendelian randomisation analysis of the effect of educational attainment and cognitive ability on smoking behaviour. *Nature Communications* **10** (2019). URL https://doi.org/10.1038/s41467-019-10679-y.

[63] LaBianca, S. *et al.* Polygenic profiles define aspects of clinical heterogeneity in attention deficit hyperactivity disorder. *Nature Genetics* **56**, 234–244 (2023). URL https://doi.org/10.1038/s41588-023-01593-7.

[64] Didelez, V. Causal concepts and graphical models. In Maathuis, M., Drton, M., Lauritzen, S. & Wainwright, M. (eds.) *Handbook of Graphical Models*, 353–376 (Chapman and Hall/CRC, Boca Raton, FL, 2018). URL https://doi.org//10.1201/9780429463976-15.

[65] Pearl, J. An introduction to causal inference. *The International Journal of Biostatistics* **6** (2010). URL https://doi.org/10.2202/1557-4679.1203.

[66] Harville, D. *Matrix Algebra from a Statistician's Perspective* (Springer-Verlag, New York, 1997). Repr. with corrections.

[67] Castelletti, F. & Consonni, G. Bayesian inference of causal effects from observational data in Gaussian graphical models. *Biometrics* **77**, 136–149 (2021). URL https://doi.org/10.1111/biom.13281.

[68] Chickering, D. M. Learning equivalence classes of bayesian-network structures. *Journal of Machine Learning Research* **2**, 445–498 (2002). URL http://www.ai.mit.edu/projects/jmlr/papers/volume2/chickering02a/chickering02a.pdf.

[69] Andersson, S. A., Madigan, D. & Perlman, M. D. A characterization of Markov equivalence classes for acyclic digraphs. *Annals of Statistics* **25**, 505–541 (1997). URL https://doi.org/10.1214/aos/1031833662.

[70] Castelletti, F., Consonni, G., Vedova, M. L. D. & Peluso, S. Learning Markov equivalence classes of Directed Acyclic Graphs: An objective Bayes approach. *Bayesian Analysis* **13**, 1235–1260 (2018). URL https://doi.org/10.1214/18-ba1101.

[71] Bottolo, L. & Richardson, S. Evolutionary stochastic search for Bayesian model exploration. *Bayesian Analysis* **5**, 583–618 (2010). URL https://doi.org/10.1214/10-ba523.

[72] Kalisch, M. & Bühlmann, P. Estimating high-dimensional Directed Acyclic Graphs with the PC-algorithm. *Journal of Machine Learning Research* 613–636 (2007). URL https://jmlr.org/papers/volume8/kalisch07a/kalisch07a.pdf.

[73] Zeng, P., Shao, Z. & Zhou, X. Statistical methods for mediation analysis in the era of high-throughput genomics: Current successes and future challenges. *Computational and Structural Biotechnology Journal* **19**, 3209–3224 (2021). URL https://doi.org/10.1016/j.csbj.2021.05.042.

[74] Zuber, V. *et al.* High-throughput multivariable Mendelian randomization analysis prioritizes apolipoprotein B as key lipid risk factor for coronary artery disease. *International Journal of Epidemiology* **50**, 893–901 (2021). URL https://doi.org/10.1093/ije/dyaa216.