**Supplemental information for**

# Understanding the assumptions underlying Mendelian Randomization

## *Table of contents*

## LD-induced pleiotropy

The scenario of LD-induces pleiotropy as depicted in Figure 2d, with variant $j$ in LD with two separate causal variants for $X$ and $Y$, is a specific violation of the instrumental variable assumptions that is comparatively easy to test.

If we have a causal variant $k$ for $X$ that is a valid instrument (and thus $\beta_k = \frac{\gamma_{Yk}}{\gamma_{Xk}} = \beta_{XY}$, then for any variant $j$ in LD with it (and with no other causal variants), we have $\gamma_{Xj} = r_{jk}\gamma_{Xk}$ and $\gamma_{Yj} = r_{jk}\gamma_{Yj}$ and therefore $\beta_j = \frac{\gamma_{Yj}}{\gamma_{Xj}} = \frac{r_{jk}\gamma_{Yj}}{r_{jk}\gamma_{Xj}} = \beta_k = \beta_{XY}$, and thus variant $j$ is also a valid instrument. Indeed, this still applies when $j$ is in LD with multiple valid instruments, in which case the $r_{jk}$ in these equations simply gets replaced everywhere with the total summed LD of $j$ with all those causal variants and similarly cancels out again in the ratio.

In case of the scenario in Figure 2d however (assuming $\beta_{YX} = 0$ to simplify notation), we have $\gamma_{Xj} = r_{jk}\alpha_{Xk}$ and $\gamma_{Yj} = r_{jk}\alpha_{Xk}\beta_{XY} + r_{ij}\alpha_{Yi}$, and therefore $\beta_j = \beta_{XY} + \frac{r_{ij}}{r_{jk}}\frac{\alpha_{Yi}}{\alpha_{Xk}}$. As such, none of the variants in the LD block are valid instruments, including variant $k$ which will be in LD with variant $i$ in this case as well.

Moreover, the $\beta_j$ will vary across variants in the LD block as a function of the relative ratio of LD $\frac{r_{ij}}{r_{jk}}$ that a variant has with the two causal variants. As such we can test the heterogeneity of the $\beta_j$ within an associated locus to determine whether this LD-induced pleiotropy is present, and can discard any loci from the MR analysis if they show evidence of this. This is what the HEIDI test in SMR[1] does (note that this is a different application from the HEIDI test in GSMR[2], where it is used as a general heterogeneity test across loci).

## Median- and mode-based methods

When using multiple variants for an MR analysis, if we assume that at least a subset of those variants are valid instruments, as noted in the main text, one approach we can take is to explicitly filter out all variants with heterogeneous $\beta_j$. We can then obtain $\beta_{XY}$ as the shared $\beta_j$ of the remaining homogeneous set of variants. However, an alternative to this approach is to make a more specific assumption about that valid subset, and to then use the median[3] or mode[4–6] of the $\beta_j$ of all variants to obtain $\beta_{XY}$ without any explicit filtering.

The median-based approach proceeds as follows. For any valid instrument we know that $\beta_j = \beta_{XY}$. If we assume that at least a majority of the variants used is a valid instrument, it must be the case that if we were to sort all their $\beta_j$ by value then the middle $\beta_j$ (if using an odd number of variants) or middle two $\beta_j$ (if using an even number of variants) must equal $\beta_{XY}$. For example, if we have 100 variants, then even if 49 of those are invalid instruments all with $\beta_j < \beta_{XY}$, ordered by value these will fill the first 49 places with the 50th and 51st being $\beta_{XY}$. Hence, the median of all the $\beta_j$ will also equal $\beta_{XY}$.

For the mode-based approaches, we can first note that we can write $\beta_j = \beta_{XY} + \delta_j$ for some deviation term $\delta_j$ for every variant $j$. We can now partition all the variants into subsets based on their $\delta_j$. All the valid instruments will define one subset with $\delta_j = 0$, whereas variants conforming to the reverse causation (Figure 1b) and mediating confounder (Figure 1c) scenarios will each define a subset as well (with separate subsets for each such mediating confounder). Other variants will generally each define their own subset containing a single variant, since it is unlikely that multiple $\beta_j$ will be the same by chance.

We can now assume that the largest of these subsets is the one with $\delta_j = 0$, ie. the most commonly occurring deviation is 0; this is also referred to as assuming that a plurality of variant are valid instruments, and specifically as ZEMPA (Zero Modal Pleiotropy Assumption) in various methods. If this assumption holds, inherently 0 is the most common value among the $\delta_j$ and therefore $\beta_{XY}$ is the most common value among the $\beta_j$, and thus the mode of the $\beta_j$ must equal $\beta_{XY}$.

## Collider bias

For a set of three variables A, B and C, if B has a direct causal relation with both A and C then there are three basic scenarios that can arise (irrespective of direct relations between A and C, and for simplicity ignoring reciprocal causation):
- If B is caused by one of the other two variables but causal for the other, then B is a mediator on the causal path between A and C
  - A -> B -> C or C -> B -> A
- If B causes both A and C, then B is a confounder of A and C
  - A <- B -> C
- If B is caused by both A and C, it is a collider for A and C
  - A -> B <- C

Collider bias, also referred to as selection bias, is a bias in the association between two variables when conditioning on a variable that is a collider for those variables. In the typical example collider bias creates an association where none exists, but it can equally amplify, reduce or remove an existing association. The same thing happens when conditioning on a descendant of a collider as well (ie. a variable upon which the collider has a causal effect), though the strength of the bias will depend on the strength of the relation between the collider and its descendant.

Note that this kind of bias can occur as a result of explicitly conditioning on a collider in a statistical analysis, but can also result from (implicit) selection mechanisms that exist as part of the data collection (eg. samples of older individuals implicitly select for longevity).

The intuition behind this kind of bias can be illustrated with a simple example. Suppose we have two independent genetic variants A and C, both capable of causing phenotype B (100% penetrance) and each with effect allele frequency of 10%. This gives the following population distribution:

|       | A = 0 | A = 1 |     |
|-------|-------|-------|-----|
| C = 0 | 81%   | 9%    | 90% |
| C = 1 | 9%    | 1%    | 10% |
|       | 90%   | 10%   |     |

In the population as a whole, there is no relation between A and C. In general, there is a 10% chance of C being 1; this is still 10% if we condition on A = 0 or A = 1, knowing the value of A gives us no information about C (and vice versa).

This changes if we select on B = 1 however (highlighted in orange; assuming no other causes of B exist). In this subpopulation, there is a 9/19 = 47.4% chance that C = 1. But now, if we know that A = 0 the chance of C = 1 becomes 100%. Given that someone has the phenotype B, something must have caused B. And if we know that it wasn't A, then it must have been C (and vice versa). Conversely, if we know that A = 1 there is no further reason to suspect that C = 1, and hence in that case the chance of C = 1 drops from 47.4% back down to the 10% level in the full population.

In this way, selecting on the collider B creates a (negative) relation between A and C that does not exist in the underlying population, and that does not reflect a causal relation between A and C (and hence can be considered spurious). In practice the relations between variables will usually be less deterministic, but the same general logic applies.


## Negative control populations

As noted in the main text, the central requirement for a negative control population is that the exposure $X$ is constrained to a particular value. This is because if $X$ is truly constrained, its value is simply held constant regardless of any causal effects trying to operate on it, effectively cancelling out those causal effects.

The causal effects of $X$ on other variables still exist however, but since the value of $X$ is the same for all individuals in the sample $X$ has no variance. It therefore also has no covariance with the variables it would causally affect, and only shifts the overall mean of those variables (assuming simple linear relations). This is the reason it does not necessarily matter what value $X$ is constrained at, unlike a control group in RCT. Whereas a control group serves as a baseline to compare the treatment groups to, the negative control population in MR only serves to evaluate the presence of pleiotropic effects of variants on $Y$ that bypass $X$.

### Constraint versus selection
Given the considerable utility of negative control populations for validating genetic instruments, it would be tempting to simply select a subpopulation where eg. $X = 0$ to use as negative control population. But because this is selected for a value rather than constrained at $X = 0$, this does not work. If $X$ is constrained to 0, this means that it is essentially fixed at that value, regardless of any causal effects operating on it. However, if it is selected to be 0, this means that either there are no causal effects moving it away from 0 to begin with (assuming $X = 0$ is in some way the 'default' state), or the causal effects operating on $X$ for these individuals happen to be balancing themselves at 0. If the latter, this means that by selecting for $X = 0$ we would implicitly also be selecting in some way for combinations of variables that causally affect $X$.

This is highly likely to result in collider bias. Suppose that we have a variant $j$ that is a valid instrument as per Figure 1a, and is a causal variant for $X$ (the same still applies if it is only in LD with the causal variant). Assume that the other assumptions in Table 1 hold, and for simplicity also that $\beta_{YX} = 0$. In this scenario, $X$ is both a mediator of the causal effect of $G_j$ on $Y$, as well as a collider for $G_j$ and $C$.

The marginal association between $G_j$ and $Y$ in this case is simply $\gamma_{Xj} = \alpha_{Xj}\beta_{XY}$, reflecting the mediated effect via $X$. If we were to now condition on $X = 0$ this mediated effect will be removed, but at the same time due to the collider bias an association $\gamma_{Cj} = -\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\beta_{CX}$ arises between $G_j$ and $C$, and therefore the conditional association with $Y$ becomes $\gamma_{Yj} = -\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\beta_{CX}\beta_{CY}$ rather than 0. As such, even for valid instruments this will lead to $\gamma_{Yj}^{(C)} \neq 0$.

### Testing validity versus estimating unmediated effects
When a likely suitable negative control sample is available, it will similarly be tempting to use an approach like PRMR[7] to estimate the unmediated part $\delta_{Yj}$ of $\gamma_{Yj}$, that is $\delta_{Yj} = \gamma_{Yj} - \gamma_{Xj}\beta_{XY}$ and use this as a correction, rather than merely testing whether variant $j$ is a valid instrument. As noted in the main text however, this type of correction is much more vulnerable to generating false positives. Whether a

population is a genuine negative control ultimately cannot be verified and must be assumed, which may itself be an invalid assumption. In this case $\gamma_{Yj}^{(C)}$ will likely not equal the $\delta_{Yj}$ in the original population.

When merely testing $\delta_{Yj} = 0$, since $\gamma_{Yj}^{(C)}$ will be biased away from 0 in this case, this will generally result in false positives on the test of $\delta_{Yj} = 0$, which in turn leads to the variant $j$ being incorrectly dismissed as an invalid instrument even if it is valid. In terms of the MR analysis itself, this will tend to translate to essentially a form of false negative, as it may result in all variants being marked as invalid instruments and therefore no test of $\beta_{XY}$ being performed.

It is of course possible that the bias is such that for an invalid instrument it exactly equals $-\delta_{Yj}$ and therefore yields a $\gamma_{Yj}^{(C)}$ of 0, in which case the variant may be incorrectly concluded to be a valid instrument. In practice it is not very likely that the unmediated component will be exactly cancelled out by such a bias however, particularly for multiple independent variants, so when used carefully and with well-powered control samples this shouldn't be an enormous risk. While it may be concluded that on the basis of the control population no valid instruments can be found, this would not lead to biased estimates of, or invalid inference on, the causal effect.

When using the negative control population to estimate $\delta_{Yj}$ however, such as in PRMR[7], no clear indicator of the control population as a whole being invalid is available. The aim is to use the estimate of $\delta_{Yj}$ to be able to use any variant regardless of whether it is a valid instrument, so there is no expectation that $\delta_{Yj}$ should be zero for any variant used. This can therefore easily lead to bias in the estimate of the causal effect.

An additional complication is that as noted, estimating $\delta_{Yj}$ requires the additional assumption that there is no reciprocal causation. The reason for this is that in the general scenario of Figure 2a $\gamma_{Yj} = \frac{1}{1-\beta_{XY}\beta_{YX}}\left(\alpha_{Yj} + \alpha_{Xj}\beta_{XY} + \alpha_{Cj}(\beta_{CY} + \beta_{CX}\beta_{XY})\right)$, which means that $\delta_{Yj} = \frac{1}{1-\beta_{XY}\beta_{YX}}\left(\alpha_{Yj} + \alpha_{Cj}\beta_{CY}\right)$. But since $\gamma_{Yj}^{(C)} = \alpha_{Yj} + \alpha_{Cj}\beta_{CY}$ as before, it does not equal $\delta_{Yj}$ anymore because the scaling factor $\frac{1}{1-\beta_{XY}\beta_{YX}}$ is missing from $\gamma_{Yj}^{(C)}$.

This scaling factor represents the rescaling of direct effects on $Y$ that occurs in this kind of scenario: any direct effect on $Y$ results in an effect on $X$ as well, which then imparts additional effect on $Y$ again (and so on; mathematically, this yields a geometric series of the causal parameters). If the signs of the causal effects $\beta_{XY}$ and $\beta_{YX}$ are the same, $\frac{1}{1-\beta_{XY}\beta_{YX}}$ will be greater than one and the initial direct effects are essentially amplified. If the signs differ $\frac{1}{1-\beta_{XY}\beta_{YX}}$ will be between 0 and 1 however, resulting in a dampening of the initial direct effects instead.

With $X$ constrained and all its causal effects blocked, this amplification/dampening disappears and hence $\gamma_{Yj}^{(C)}$ only reflects the initial, direct effects on $Y$. The reason this is not a problem when just testing $\delta_{Yj} = 0$ is that under that null hypothesis there is no direct effect to be amplified or dampen, and hence the scaling factor is irrelevant in that case.

### Using gene-environment interactions
It has been suggested that gene-environment interaction can be used to similar effect as methods like PRMR[7], but without the need for a negative control population[8]. The intuition behind this is that if there is a variable $Z$ that modifies the effect of variant $j$ on the exposure in the form of a linear interaction, then this creates a so-called 'no relevance group', a subgroup within the population where the instrument and exposure are independent. The association in that subgroup with the outcome could then be used to estimate the unmediated component $\delta_{Yj}$. If this subgroup does not actually exist because it falls outside the range of $Z$ within the population, then in that case it would still be possible to extrapolate to define a hypothetical 'no relevance group' (assuming the interaction is indeed fully linear).

In practice however, this approach requires a number of assumptions that PRMR does not, and effectively just reduces to an instrumental variable analysis using the interaction term $I_j = G_j Z$ of $Z$ and variant $j$ as the instrument rather than $G_j$ itself. It therefore requires all the same assumptions to be satisfied for $I_j$ as a regular MR analysis does for $G_j$ (plus additional assumptions pertaining to the interaction), and runs into the same issues when these assumptions fail. For example, in just the same way as $G_j$ itself, $I_j$ may be directly associated with a confounder $C$ rather than with $X$, which would be similarly difficult to detect.

To better understand such a scenario, we can simplify somewhat by supposing that the variable $Z$ is a discrete and bounded integer, and thus divides the population in an ordered set of subpopulations for each value $z$ that $Z$ can take. Within each of these subpopulations we will further assume the causal effects between $X$, $Y$ and $C$ to be the same, though in practice these may well interact with $Z$ themselves as well.

In the case that $I_j$ is indeed a valid instrument, we would have subgroup-specific effects $\alpha_{Xj}^{(z)} = z\alpha_{IXj}$, where $\alpha_{IXj}$ is the effect of $I_j$ on $X$ (since for $Z = z$, we have $I_j\alpha_{IXj} = G_j z\alpha_{IXj} = G_j\alpha_{Xj}^{(z)}$). For further simplicity we will assume that there is no main effect of $G_j$ on $X$, and as such the $Z = 0$ group will be identified as our 'no relevance group'. For the observable marginal effects, we therefore have $\gamma_{Xj}^{(z)} = \alpha_{Xj}^{(z)}$ and $\gamma_{Yj}^{(z)} = \alpha_{Xj}^{(z)}\beta_{XY} = \gamma_{Xj}^{(z)}\beta_{XY}$, and hence indeed $\frac{\gamma_{Yj}^{(z)}}{\gamma_{Xj}^{(z)}} = \beta_{XY}$ (except for $Z = 0$, where there this would yield a division by 0).

We can similarly see how this would be able to resolve a scenario like that in Figure 2d as well, where there are additional direct associations between $G_j$ and $Y$. In this case $\gamma_{Yj}^{(z)} = \gamma_{Xj}^{(z)}\beta_{XY} + \alpha_{Yj}$, and hence in our 'no relevance group' $\gamma_{Yj}^{(0)} = \alpha_{Yj}$. It therefore follows that for the other values of $Z$, $\frac{\gamma_{Yj}^{(z)} - \gamma_{Yj}^{(0)}}{\gamma_{Xj}^{(z)}} = \beta_{XY}$, in much the same way that PRMR uses the negative control population to remove bias.

However, from the expression $\gamma_{Yj}^{(z)} = \gamma_{Xj}^{(z)}\beta_{XY} + \alpha_{Yj}$ it is also easy to spot a problem: we must assume that $\alpha_{Yj}$ is in fact constant across values of $Z$, which is by no means guaranteed to be the case. And one particular way in which this may occur is the mediating confounder scenario in Figure 1c (where for simplicity again we will assume $\beta_{YX} = 0$). In this case the interaction affects $C$ rather than $X$ directly, such that $\alpha_{Cj}^{(z)} = z\alpha_{ICj}$. This gives us $\gamma_{Xj}^{(z)} = \alpha_{Cj}^{(z)}\beta_{CX}$ and $\gamma_{Yj}^{(z)} = \alpha_{Cj}^{(z)}(\beta_{CX}\beta_{XY} + \beta_{CY})$. This would yield $\gamma_{Yj}^{(0)} = 0$ and hence $\frac{\gamma_{Yj}^{(z)} - \gamma_{Yj}^{(0)}}{\gamma_{Xj}^{(z)}} = \beta_{XY} + \frac{\beta_{CY}}{\beta_{CX}}$, exactly the same biased value that we run into with this scenario in regular MR.

Conceptually speaking, the issue that we run into here is similar to the problem that arises in the context of using an $X = 0$ subgroup of a population as a negative control, as discussed above (although the math works out differently here). Both approaches inherently use a form of selection, whereas the inferential strength of negative control populations and indeed of control groups in experimental designs comes from the imposition of constraints.

When selecting a group for having a particular property (like $X = 0$ or $\gamma_{Xj} = 0$), this only means that the balance of relevant forces contrived to have this particular group have that particular (parameter) value. A genuine constraint however, simply overrides those normally relevant forces and imposes the particular value on that group.

Statistically we could still formulate the resulting situation as an interaction as we did now, by using a population indicator variable to define $Z$ anddesignating the negative control population as the $Z = 0$ group with $\gamma_{Xj}^{(0)}$ similarly 0, and this may give the appearance of similarity. Yet by the same token we can use the same ANOVA model to compare experimental groups as we do with observational data.
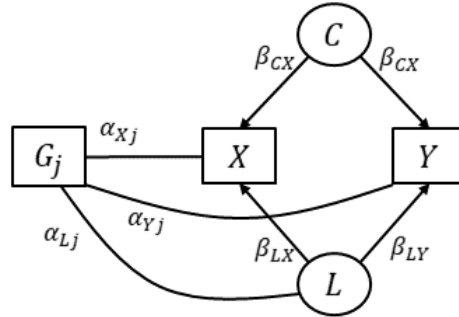
There we use the knowledge, not captured by the statistical model, that those experimental conditions were actively imposed. This is the same principle that gives negative control populations their inferential strength. And extending the comparison, it is limited in its inferential strength to the extent that we do not know whether there are other (relevant) structural differences between the populations, in the same way as experimental designs are limited by any factors that were not fully controlled across the conditions.

## Related non-MR methods

The strong reliance on assumptions to draw causal inference is not a feature unique to MR, and affects other methods developed for such purposes as well. For these, we can similarly examine the underlying model assumed by the method to determine what role the assumptions made play in the inference and how the method behaves under alternate scenarios. We do so here for two methods also intended for causal inference based on genetic data. As in the main text we do this under the idealized scenario that associations between observed variables are fully known, and unless stated otherwise will also assume that all the assumptions in Table 1 other than the instrumental variable assumptions all hold.

### *Latent Causal Variable*
The Latent Causal Variable (LCV) model[9] is a whole-genome model presented as an alternative to MR. The same general model is assumed for every variant $j$, with effects of $G_j$ directly on $X$ and $Y$ as well as mediated by a latent variable $L$, as per the graph below, with no direct effects modelled between $X$ and $Y$. The effects of the variants are modeled as random, and it is assumed that $\alpha_{Lj}$ is independent of $\alpha_{Xj}$ and $\alpha_{Yj}$ (note that for the estimation, the marginal effects $\gamma_X$ and $\gamma_Y$ are also assumed to be standardized, as is the distribution of $\alpha_{Lj}$).



For a given variant we can see that for a variant $j$, the (unstandardized) marginal effects are $\gamma_{Xj} = \alpha_{Lj}\beta_{LX} + \alpha_{Xj}$ and $\gamma_{Yj} = \alpha_{Lj}\beta_{LY} + \alpha_{Yj}$. When assuming the different $\alpha_j$ parameters are independent of each other, it follows that $\text{cov}(\gamma_{Xj}, \gamma_{Yj}) = \beta_{LX}\beta_{LY}$ (and similiarly $\text{var}(\gamma_{Xj}) = \beta_{LX}^2 + \text{var}(\alpha_{Xj})$ and $\text{var}(\gamma_{Yj}) = \beta_{LY}^2 + \text{var}(\alpha_{Yj})$). Effectively, the correlation of the marginal associations of $G_j$ with $X$ and $Y$, strongly related to the genetic correlation of $X$ and $Y$, is decomposed into a shared and unique part, with the shared part attributed to the latent variable $L$ (although the above equations show only three known values for four unknown parameters, this issue is resolved by LCV by using higher order product moments per variant as additional statistics).

The main parameters of interest are $\beta_{LX}$ and $\beta_{LY}$, and in particular the genetic causality proportion (GCP) statistic $\text{GCP} = \frac{\log|\beta_{LX}| - \log|\beta_{LY}|}{\log|\beta_{LX}| + \log|\beta_{LY}|}$. This reflects the relative size of $\beta_{LX}$ and $\beta_{LY}$ on a scale of -1 to 1,

being 0 if they are equal in size, going towards 1 if $\beta_{LX}$ is larger and towards -1 if $\beta_{LY}$ is larger. The authors propose an interpretation of the GCP statistic in terms of 'genetic causality', suggesting that $X$ is partially genetically causal for $Y$ if $L$ has a stronger genetic correlation with $X$ than with $Y$, ie. $\text{GCP} > 0$ (and vice versa).

This interpretation does not follow from the model used, however. In a literal sense, the model even excludes the existence of causal effects between $X$ and $Y$ a priori. The apparent reasoning is that if such causal effects do exist, in the model these will instead essentially be absorbed into $L$, with causal effects of $X$ on $Y$ effectively pulling $L$ 'closer' to $X$ (ie. greater $\beta_{LX}$ and smaller $\beta_{LY}$) and effects of $Y$ on $X$ doing the reverse. This is also reflected in the use of the GCP statistic, which reflects the relative skew of $\beta_{LX}$ and $\beta_{LY}$ on a standardized scale.
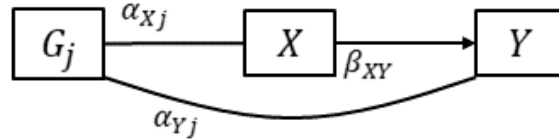
In the case that all shared effects are indeed mediated by $X$ on $Y$ (ie. all variants associated with $X$ conform to the model in Figure 1d), $L$ is effectively absorbed into $X$ and the GCP becomes 1. Equivalently, if all effects are mediated in the opposite direction (with all variants associated with $Y$ conforming to Figure 1e), the GCP becomes -1. As such, if the GCP is at either extreme, this can indeed be taken as evidence for a causal effect of $X$ on $Y$ or vice versa.

In general however the $\beta_{LX}$ and $\beta_{LY}$ parameters will reflect some average of the effects across all the heritable confounders of $X$ and $Y$, as well as a mixture of direct causal effects between $X$ and $Y$ (if present). The GCP then only reflects the asymmetry in the effects of $L$ on $X$ and $Y$, but provides no further basis for meaningful causal inference. And although to an extent the asymmetry itself may still be somewhat informative, this can also induced by imperfect observation of the causally relevant instances of $X$ and $Y$ as also discussed in the main text, and may therefore just be due to for example greater measurement error for one of the variables.


***Genetic Instrumental Variable regression***

In the idealized version of the Genetic Instrumental Variable regression (GIV) model[10], a regression $Y = \theta_X X + \theta_R R_{Y|X} + \varepsilon$ is performed, where $R_{Y|X}$ represents the total genetic effect on $Y$ correcting for $X$. The estimate of $\theta_X$ is then used as an estimate of $\beta_{XY}$. The idea behind this approach is that the conditional genetic effect $R_{Y|X}$ would capture all genetic associations with $Y$ that are not mediated by $X$. When this term is then included when regressing $Y$ on $X$, these pleiotropic associations are corrected for. Effectively, $R_{Y|X}$ is intended to capture all confounding of $X$ and $Y$ by the genetic effects of variants that causally affect both phenotypes.

This would work if the causal model for all variants corresponds to the graph below.



In the presence of confounders however, even if these are not heritable, this would run into the problem of collider bias. Assume that all variants associated with $X$ are valid instruments corresponding to Figure 1d, and simplify the mathematics by also assuming there is only a single confounder $C$ and no LD between the variants. In this case, for a variant $j$ the association of $G_j$ with $Y$ given $X$ would be $-\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\beta_{CX}\beta_{CY}$.

We therefore obtain $R_{Y|X} = -\beta_{CX}\beta_{CY}\sum_j \frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}G_j$, with $\mathrm{Var}(R_{Y|X}) = \beta_{CX}^2\beta_{CY}^2\sum_j\left(\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\right)^2$, $\mathrm{Cov}(X,R_{Y|X}) = -\beta_{CX}\beta_{CY}\sum_j\frac{\alpha_{Xj}^2}{1-\alpha_{Xj}^2}$ and $\mathrm{Cov}(Y,R_{Y|X}) = -\beta_{CX}\beta_{CY}\beta_{XY}\sum_j\frac{\alpha_{Xj}^2}{1-\alpha_{Xj}^2}$. For the regression given above it therefore follows that

$$\theta_X = \frac{(\beta_{XY}+\beta_{CX}\beta_{CY})\beta_{CX}^2\beta_{CY}^2\sum_j\left(\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\right)^2 - \beta_{XY}\left(\beta_{CX}\beta_{CY}\sum_j\frac{\alpha_{Xj}^2}{1-\alpha_{Xj}^2}\right)^2}{\beta_{CX}^2\beta_{CY}^2\sum_j\left(\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\right)^2 - \left(\beta_{CX}\beta_{CY}\sum_j\frac{\alpha_{Xj}^2}{1-\alpha_{Xj}^2}\right)^2} = \beta_{XY} +$$

$$\frac{\beta_{CX}^3\beta_{CY}^3\sum_j\left(\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\right)^2}{\beta_{CX}^2\beta_{CY}^2\sum_j\left(\frac{\alpha_{Xj}}{1-\alpha_{Xj}^2}\right)^2 - \left(\beta_{CX}\beta_{CY}\sum_j\frac{\alpha_{Xj}^2}{1-\alpha_{Xj}^2}\right)^2} = \beta_{XY} + \frac{\beta_{CX}\beta_{CY}\mathrm{Var}(R_{Y|X})}{\mathrm{Var}(R_{Y|X})-\mathrm{Cov}(X,R_{Y|X})^2} = \beta_{XY} + \beta_{CX}\beta_{CY}\frac{1}{1-\mathrm{Cor}(X,R_{Y|X})^2}.$$

In other words, the slope of the regression would be biased by a term $\beta_{CX}\beta_{CY}\frac{1}{1-\mathrm{Cor}(X,R_{Y|X})^2}$ even if all variants are valid instruments, unless no confounding is present for $X$ and $Y$ at all, which is exceedingly unlikely to be the case for any $X$ and $Y$.

Note also that in practice $R_{Y|X}$ would not be available in practice, it is instead constructed as a PRS $\hat{R}_{Y|X}$ for $Y$, using weights from a GWAS for $Y$ that includes $X$ as a covariate. Because this would result in bias due to noise in $\hat{R}_{Y|X}$ relative to $R_{Y|X}$, an instrumental variable approach using instruments for $\hat{R}_{Y|X}$ is used in GIV to try to remove that bias.

## Supplemental references

1. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* **48**, 481–487 (2016).
2. Zhu, Z. *et al.* Causal associations between risk factors and common diseases inferred from GWAS summary data. *Nat. Commun.* **9**, (2018).
3. Bowden, J., Davey Smith, G., Haycock, P. C. & Burgess, S. Consistent Estimation in Mendelian Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genet. Epidemiol.* **40**, 304–314 (2016).
4. Qi, G. & Chatterjee, N. Mendelian randomization analysis using mixture models for robust and efficient estimation of causal effects. *Nat. Commun.* **10**, 1–10 (2019).
5. Hartwig, F. P., Smith, G. D. & Bowden, J. Robust inference in summary data Mendelian randomization via the zero modal pleiotropy assumption. *Int. J. Epidemiol.* **46**, 1985–1998 (2017).
6. Burgess, S., Zuber, V., Gkatzionis, A. & Foley, C. N. Modal-based estimation via heterogeneitypenalized weighting: Model averaging for consistent and efficient estimation in Mendelian randomization when a plurality of candidate instruments are valid. *Int. J. Epidemiol.* **47**, 1242–1254 (2018).
7. Van Kippersluis, H. & Rietveld, C. A. Pleiotropy-robust Mendelian randomization. *Int. J. Epidemiol.* **47**, 1279–1288 (2018).
8. Spiller, W., Slichter, D., Bowden, J. & Davey Smith, G. Detecting and correcting for bias in Mendelian randomization analyses using Gene-by-Environment interactions. *Int. J. Epidemiol.* **48**, 702–712 (2018).
9. O'Connor, L. J. & Price, A. L. Distinguishing genetic correlation from causation across 52 diseases and complex traits. *Nat. Genet.* **50**, 1728–1734 (2018).
10. DiPrete, T. A., Burik, C. A. P. & Koellinger, P. D. Genetic instrumental variable regression: Explaining socioeconomic and health outcomes in nonexperimental data. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E4970–E4979 (2018).