- Leveraging family data to design Mendelian Randomization that is

 provably robust to population stratification
- Nathan LaPierre^{1*†}, Boyang Fu^{1*}, Steven Turnbull², Eleazar Eskin^{1,3,4}, and Sriram Sankararaman^{1,3,4†}
- 1. Department of Computer Science; 2. Department of Statistics; 3. Department of
- 6 Computational Medicine; and 4. Department of Human Genetics, UCLA, Los Angeles CA
- * These authors contributed equally to this work.
 - † Email corresponding authors at nathanl2012@gmail.com and sriram@cs.ucla.edu

9 Abstract

Mendelian Randomization (MR) has emerged as a powerful approach to leverage genetic instruments to infer causality between pairs of traits in observational studies. However, the results of such studies are susceptible to biases due to weak instruments as well as the confounding effects of population stratification and horizontal pleiotropy. Here, we show that family data can be leveraged to design MR tests that are provably robust to confounding from population stratification, assortative mating, and dynastic effects. We demonstrate in simulations that our approach, MR-Twin, is robust to confounding from population stratification and is not affected by weak instrument bias, while standard MR methods yield inflated false positive rates. We then conduct an exploratory analysis of MR-Twin and other MR methods applied to 121 trait pairs in the UK Biobank dataset. Our results suggest that confounding from population stratification can lead to false positives for existing MR methods, while MR-Twin is immune to this type of confounding, and that MR-Twin can help assess whether traditional approaches may be inflated due to confounding from population stratification.

3 Introduction

Mendelian Randomization (MR) is a widely-used analytical tool that uses genetic variants ("genetic instruments") to determine whether one trait (the "exposure") has a causal effect on another (the 25 "outcome"). With the availability of massive biobank datasets such as the UK Biobank (Bycroft 26 et al., 2018), MR analyses have become increasingly powerful and have been used to identify causal relationships between numerous pairs of traits (Hemani et al., 2018; Wade et al., 2018; 28 Lyall et al., 2017; Haycock et al., 2017; Haase et al., 2012). The validity of MR rests on three key assumptions (Lawlor et al., 2008): (i) that the genetic instrument is significantly associated with the exposure; (ii) that the genetic instrument is independent of confounders of the exposure-31 outcome relationship; (iii) that the genetic instrument affects the outcome only through its effect 32 on the exposure. 33 Unfortunately, the latter two assumptions are often violated in practice, due to several factors 34 including horizontal pleiotropy, population stratification (and related phenomena such as assortative 35 mating and dynastic effects), and batch effects. Even when these assumptions are met, the weak effects of typical genetic instruments on the exposure coupled with spurious correlation between genetic instruments and confounders (Burgess et al., 2011) can bias the results of MR analyses ("weak instrument bias"). The problem of population stratification has been extensively studied in the Genome-Wide Association Study (GWAS) literature, and approaches for mitigating its effects have been developed, including the usage of Principal Components Analysis (PCA) and Linear Mixed Models (LMMs) (Price et al., 2010). These approaches have generally been found to be effective at reducing the confounding introduced by population stratification (Price et al., 2010). However, recent studies have demonstrated that, with sample sizes as large as those found in 44 modern biobanks, even a small amount of residual population stratification can cause a considerable amount of bias (Cook et al., 2020; Brumpton et al., 2020; Berg et al., 2019; Haworth et al., 2019), and may even cause false positives in MR analysis (Haworth et al., 2019; Cinelli et al., 2022). In addition, while the confounding effects of population stratification are well-known, less attention has been directed towards confounding from other phenomena such as (cross-trait) assortative mating and dynastic effects, which can also cause MR false positives (Brumpton et al., 2020; Hartwig et al., 2018). Recent work has demonstrated that cross-trait assortative mating is widespread and substantially inflates genetic correlation estimates between many trait pairs (Border et al., 2022).

It has recently been proposed that family-based genetic datasets could be used in MR studies to avoid confounding from population stratification (Brumpton et al., 2020; Pingault et al., 2018).

A recent suite of methods have been developed for this purpose, and were shown to reduce the bias from this type of confounding (Brumpton et al., 2020). However, like other MR methods, these methods are susceptible to weak instrument bias, which can be substantial for small family-based datasets (Brumpton et al., 2020). In this paper, we introduce MR-Twin, a test for causal effects between pairs of traits that is able to leverage family-based genetic data to provably control for population stratification and utilize publicly-available summary statistics estimated in large biobank datasets to achieve power competitive with top existing methods for the same sample size.

We develop versions of MR-Twin for trio, parent-child duo, and sibling data, evaluate MR-Twin's ability to control false positives due to population stratification and weak instrument bias, and compare it with existing methods.

65 Results

66 Methods overview

In a Mendelian Randomization (MR) analysis, we wish to determine whether one phenotype (the "exposure") has a causal effect on another phenotype (the "outcome") using genetic instrumental variables, which can be either single nucleotide polymorphims (SNPs), a polygenic score, or other genetic features. Under the assumption that the genetic instruments are associated with the expo-70 sure and are independent of the outcome given the exposure, the MR effect estimate of the exposure 71 on the outcome will be valid even if there are unobserved confounders of the exposure-outcome re-72 lationship. The independence assumption, however, is often violated by population stratification (Figure 1A) or horizontal pleiotropy, as these phenomena cause the genetic instruments to be correlated with the outcome through pathways other than those through the exposure. 75 MR-Twin is a method that uses family-based genetic data to construct a test for whether the 76 exposure has an effect on the outcome that is immune to confounding from population structure. It is based on the key idea that the genotypes of observed individuals are independent of popu-78 lation structure given the genotypes of the individuals' parents, since the mechanisms by which

genetic information is passed from parents to offspring are known (Figure 1B). In other words, conditioned on the parental genotypes, population structure provides no additional information about the distribution of the offspring's genotypes. Thus, by conditioning on the parental genotypes, 82 confounding from population stratification can be avoided (Figure 1C), along with confounding from other phenomena such as cross-trait assortative mating and dynastic effects which operate through the parental genotypes (see Figure 1 of (Brumpton et al., 2020)). 85 We now outline the algorithm in the context of a trio design in which we have genetic data on 86 the parents and the offspring. Let \mathbf{X} and O denote the genotypes and outcome phenotype values 87 respectively for some individual, and let $(\mathbf{X}_n; O_n)_{n=1}^N$, denote these across N trios. Also let **P1** and 88 P2 denote the genotypes of the parents of the individual with genotypes X, and let A := (P1, P2)89 refer to the set of parental genotypes. Let Z denote the set of external confounders measured on the same individual, which we define as the set of confounders that satisfy $X \perp Z \mid A$. Thus, population stratification is an external confounder (as are assortative mating and dynastic effects) 92 while horizontal pleiotropy is not. The key idea is that we can formulate a hypothesis test of a 93 causal effect conditional on the parental haplotypes A. Bates et al. (Bates et al., 2020) show that such a test is also a test of the stronger null hypothesis of a causal effect conditional on (A, Z). 95 The way that this is accomplished is through a conditional randomization test, similar to the 96 Digital Twin Test proposed by Bates et al in the context of GWAS (Bates et al., 2020; Candes et al., 2018). The idea is to sample so-called "digital twins" $\tilde{\mathbf{X}}$ from each set of parents \mathbf{A} such that $\tilde{\mathbf{X}} \mid \mathbf{A}$ 98 has the same distribution as $X \mid A$, which can easily be accomplished using the laws of Mendelian 99 inheritance (Methods). We construct B such random samples across all trios, $(\tilde{\mathbf{X}}_n, O_n)_{n=1}^N$, and for 100 each set b of twins we compute a test statistic $t_b = t((\tilde{\boldsymbol{X}}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}})$ representing the strength of 101 association between the genetically-predicted exposure and the outcome. We also compute a test 102 statistic for the true offspring of the trios, $t^* = t((\boldsymbol{X}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}}).$ 103 We can then obtain a p-value for a non-zero causal effect of the exposure on the outcome, 104 $p = \frac{1+1\{t_b \ge t^*\}}{1+B}$. The set of B statistics derived from the digital twins represents a null distribution 105 conditioned on the parental genotypes. If there is a true nonzero effect of the exposure on the 106 outcome, we expect the statistic derived from the true offspring to be stronger than statistics derived from digital twins whose genotypes are randomly sampled from the parental genotypes. 108 The test statistic and algorithm are explained in more detail in the Methods section. 109

MR-Twin controls for arbitrarily strong population stratification confounding in simulations 111

Algorithm 1 Simulate genotypes under population structure

```
1: procedure SIMGENO(F_{ST}, groups, N,M) \triangleright F_{ST} is the fixation index, groups is the number of
    populations, N is number of samples, M is number of SNPs.
         Initialize the average MAF \bar{\mathbf{f}}_i i.i.d. Unif(0.05, 0.5) for each SNP j.
2:
3:
         for k \le \text{groups } \mathbf{do}
              \mathbf{f}^k \sim Beta(\mathbf{\bar{f}} \frac{(1-F_{ST})}{F_{ST}}, (1-\mathbf{\bar{f}}) \frac{(1-F_{ST})}{F_{ST}})
Generate genotype matrix \mathbf{X}^{(k)} of population k such that x_{ij}^k \sim Bin(2, \mathbf{f}_j^k) for each
4:
    individual i and SNP j.
         end for
6:
         \boldsymbol{X} = [\boldsymbol{X}^{(1)}; \dots; \boldsymbol{X}^{(groups)}]
7:
                                                                                     > Stack the rows of each genotype matrix
```

Algorithm 2 Simulate population-stratified phenotypes

```
1: procedure GETPHENO(\mathbf{X}, U, h^2, \alpha_E, \gamma_{ue}, \gamma_{uo}) \triangleright \mathbf{X} is the normalized genotype matrix, U is
   a vector with the fixed population label for each sample, h^2 is heritability of exposure E, \alpha_E is
   effect of E on outcome O, \gamma_{ue} and \gamma_{uo} are fixed confounding effects of U on E and O.
```

- Generate genetic coefficient $\boldsymbol{\beta} \sim \mathcal{N}(0, h^2 \mathbf{I}_M)$ 2:
- 3:

return Genotype matrix X

- Compute $\sigma_{\epsilon_e}^2 = 1 h^2$, $\sigma_{\epsilon_o}^2 = 1 \alpha_E^2$ Simulate $E = \mathbf{X}\boldsymbol{\beta} + \gamma_{ue}U + \epsilon_e$ where $\epsilon_e \sim \mathcal{N}(0, \sigma_{\epsilon_e}^2 \mathbf{I})$ 4:
- Simulate $O = \alpha_E E + \gamma_{uo} U + \epsilon_o$ where $\epsilon_o \sim \mathcal{N}(0, \sigma_{\epsilon_o}^2 I)$ 5:
- return (E,O)

9: end procedure

8:

7: end procedure

We compared the performance of MR-Twin to other MR methods via simulations consisting 112 of two populations with allele frequency differences modeled according to the standard Balding-113 Nichols model (Balding and Nichols, 1995), following previous works (Ochoa and Storey, 2021; 114 Conomos et al., 2016; Chen et al., 2015; Hubisz et al., 2009; Price et al., 2006; Pritchard et al., 115 2000). The procedure for simulating the genotypes is outlined in Algorithm 1. We use this algorithm 116 to simulate "external" samples (non-trio data, e.g. from a biobank), as well as the parents for the 117 trios. The offspring genotypes for the trios can then be easily sampled given the parental genotypes 118 (Methods). For each sample, we retain the population label, a binary variable indicating which 119 population each sample belongs to. Unless otherwise specified, each simulation had 50,000 (false 120 positive rate simulations) or 100,000 (power simulations) external samples and 1,000 trio samples 121 evenly split between two populations with fixation index $F_{ST} = 0.01$ and 100 SNPs, 50 of which were causal for the exposure trait. Unless otherwise specified, the heritability of the exposure trait was set to $h^2 = 0.2$.

Next, we simulate both the exposure and outcome phenotypes following a linear model, as outlined in Algorithm 2. This model allows the population labels from the first step to have an effect on the exposure and outcome phenotypes, which models population stratification that violates the MR assumptions. We use this setting to assess the false positive rates (FPR) of methods under population stratification, allowing the effects of the population labels on the exposure and outcome phenotypes to range from 0 (no confounding) to 0.8 (substantial confounding). In a separate set of simulations to assess power, we set the confounding effect to 0 and varied the causal effect.

We performed 1000 simulation replicates under these settings, each time simulating a set of 132 external and trio genotypes and phenotypes according to the chosen parameters, performing linear 133 regression between each SNP and the exposure and outcome phenotypes, and using the resulting 134 association statistics as input to each of the MR methods. We excluded SNPs with association 135 p-values of above 0.05/M (M=100) with the exposure phenotypes in the external data in order to 136 limit weak instrument bias (Burgess et al., 2011). The methods we assessed include the trio mode 137 of MR-Twin, standard inverse-variance weighted (IVW) MR (Burgess et al., 2013), MR-Egger 138 (Bowden et al., 2015), the Weighted Median Estimator (Bowden et al., 2016), the Mode-Based 139 Estimator (Hartwig et al., 2017), and a method introduced by Brumpton et al (Brumpton et al., 140 2020) to use family data to control for confounding due to population stratification and other 141 population-related effects. Brumpton et al provide a suite of methods for different family datasets, 142 following previous work such as (Fulker et al., 1999); here we focus on the trio-based method they 143 describe (Brumpton et al., 2020), and simply refer to that method as "Brumpton" below.

The trio mode of MR-Twin maintained a calibrated FPR irrespective of the strength of confounding (Figure 2A). Non-family-based methods such as IVW, Egger, Median, and Mode all displayed substantially inflated FPR in the face of confounding consistent with their sensitivity to potential residual population stratification. The Brumpton method also displayed slightly inflated FPR, which increased with the strength of the confounding effect, likely due to weak instrument bias (Brumpton et al., 2020). To mitigate the impact of weak instrument bias, we applied a common approach employed in MR studies (Burgess et al., 2011) that involves filtering out variants for which the F-statistic of the association signal is low (F < 10 following previous recommendations). This rendered the FPR inflation negligible, but also rendered Brumpton substantially less powerful than MR-Twin, whereas the "unfiltered" mode had similar power to MR-Twin (Figure 2B). Results with confidence intervals are shown in Figure S10. We further investigated the weak instrument bias by running simulations with no SNP filtering based on external data and with increasing numbers of SNPs – settings expected to generate large numbers of weak instruments – and confirmed that Brumpton had greater FPR inflation in these settings while MR-Twin remained calibrated (Figure S11) and did not lose power (Figure S12).

The standard MR methods (IVW, Egger, Median, and Mode), when run on the external data, 160 had substantially higher power than the family-based methods, MR-Twin and Brumpton (Figure 161 2B). We performed additional simulations to understand if the lower power of MR-Twin was due to 162 the smaller number of trios as opposed to methodological limitations. When applied to the offspring 163 in each trio (Figure 3), the standard MR methods still had substantially inflated FPR (Figure 3A) but similar power to MR-Twin and Brumpton (Figure 3B). We also evaluated the FPR and power 165 of these methods under varying number of trios (Figure S1). We observed that increasing number 166 of trios increased power for all methods, as expected, suggesting that the family-based methods can 167 be expected to obtain increased power as more genetic family data are ascertained in the future. 168 The relative power of the methods remained roughly consistent across these experiments. 169

We also evaluated the Area Under the Receiver Operating Characteristic Curve (AUC-ROC: 170 Figures S5 and S6). Comparing the two main family-based approaches, MR-Twin generally had 171 higher AUC than Brumpton (filtered). Predictably, the AUC of standard MR methods drops 172 sharply when there is confounding, and MR-Twin outperforms these methods in most such cases, 173 though Egger was a notable exception in our findings and remained competitive with family-based approaches even under confounded settings. Similarly to our findings in Figures 2 and 3, family-175 based methods are more competitive with standard MR methods when run on similar sample sizes. 176 As an additional sensitivity analysis, we also assessed the FPR (Figure S7) and Power (Figure S8) 177 of methods in settings where there are very few instruments or very low heritability. The trends 178 were broadly similar to those seen in Figures 2 and 3. MR-Twin maintained a calibrated FPR in 170 all settings, although it did suffer a loss of power when heritability was very low (Figure S8B and 180 S8D). 181

We performed simulations increasing the magnitude of population structure as measured by

182

the F_{ST} (without necessarily increasing the confounding strength), and observed that increasing the population structure leads to further FPR inflation for standard MR methods (Figure S2). We observed inflated FPR for standard MR methods even when there is no confounding (stratification) for large values of F_{ST} (Figure S2) likely due to correlation or linkage disequilibrium among the genetic variants induced by population structure. The standard implementation of IVW and Egger (Yavorska and Burgess, 2017; Broadbent et al., 2020) allows the user to pass in a variant correlation matrix, which removed the FPR inflation with no stratification (Figure S2); other methods such as Median and Mode do not currently have this option.

Next we assessed the runtime of methods run on the trio data (Figure S3). Brumpton, along 191 with non-trio based methods (e.g., IVW), had similar run times (<1-5 seconds per simulation 192 replicate); for succinctness, only Brumpton is shown. MR-Twin (with 100 simulated digital twins) 193 took roughly one minute per simulation replicate under the simulation settings described above, 194 with time increasing to up to four minutes if the number of families or SNPs was increased. The 195 number of digital twins to simulate for MR-Twin involves a trade-off between speed and stability of 196 results. We assessed the stability of MR-Twin with different numbers of digital twins, with results 197 shown in Figure S9. We interpreted these findings as indicating that 100 digital twins are likely 198 stable enough for simulations for which many replicates are run and speed is a priority, but 1000 199 or more digital twins is recommended for one-off real data analysis. Therefore we simulated 100 200 digital twins in our simulations and 1000 in our real data analysis. We note, however, that while the 201 MR-Twin runtime increases linearly with the number of digital twins simulated, the generation and 202 statistic computation for the digital twins can be done in parallel, so many twins can be simulated 203 efficiently given multiple compute cores or nodes. For clarity of results, we did not take advantage 204 of this in our runtime assessment. 205

Finally, MR-Twin also enables users to use parent-child duo or sibling datasets (Methods). We assessed the performance of these modes versus the trio mode of MR-Twin (Figure S4). We found that the duo and sibling modes, while having lower FPR than most standard MR methods, did not maintain a calibrated FPR at high levels of confounding, which is expected since the precise sampling of offspring genotypes from parents is not possible when either or both of the parental genotypes are not available.

206

207

208

200

210

12 Application to trio data in the UK Biobank

In order to assess the results given by MR-Twin relative to other approaches in a real data context, 213 we next applied MR-Twin and four other MR methods (IVW, Egger, Median, and the Brumpton et 214 al. method (Brumpton et al., 2020)) to 144 real trait pairs in the UK Biobank (Bycroft et al., 2018). 215 These consisted of all pairwise combinations of 12 metabolic, anthropometric, and socioeconomic 216 traits that were widely measured among the UK Biobank participants (listed in Table S1). We 217 isolated 955 White British genetic trios from the full UK Biobank dataset (Supplemental Materials) 218 and used PLINK (Purcell et al., 2007) to run linear regression on the remaining unrelated White 219 British individuals for these 12 traits, including the top 20 principal components (PCs), age, and 220 sex as covariates. The genetic instruments selected for each analysis were the SNPs with genome-221 wide significant p-values ($< 5.0 \times 10^{-8}$) for the exposure trait, after linkage disequilibrium (LD) 222 pruning was performed so that none of these instruments were in substantial LD with one another 223 (Supplemental Materials). Ignoring the degenerate cases where the exposure and outcome were the same trait or where there were no significant SNPs for the exposure trait (as was the case for the 225 Townsend Deprivation Index [TDI]), there were 121 usable trait pairs. 226 Table 1 shows the results for six selected trait pairs (excluding Median for brevity because it gave 227 similar results to IVW), while Supplemental Table S2 shows the full set of analyses. (Brumpton 228 was run with several different variant filtering settings to assess the impact of potential weak 220 instrument bias (Supplemental Materials); results for all runs are given in Supplemental Table S2.) 230 For Table 1, we selected six analyses: two positive controls representing causal effects that are true 231 by definition (LDL Cholesterol \rightarrow Total Cholesterol and Weight \rightarrow Body Mass Index [BMI]), two 232 negative controls that represent seemingly implausible effects (Glucose \rightarrow TDI and Height \rightarrow Body 233 Fat), and two trait pairs with unclear or conflicting evidence (BMI → Diastolic Blood Pressure [DBP] and $BMI \rightarrow TDI$). In particular, previous studies have identified a significant effect for BMI 235 \rightarrow DBP (Lyall et al., 2017) and for BMI \rightarrow TDI in women (Tyrrell et al., 2016) with IVW analysis, 236 although Egger analysis did not replicate the significant findings in either case (Lyall et al., 2017; Tyrrell et al., 2016). 238 All methods performed as expected on the controls, with highly significant p-values for positive 239

controls and insignificant p-values for negative controls. For BMI \rightarrow DBP, IVW and Brumpton

yielded significant results while Egger and MR-Twin did not. For BMI → TDI, IVW and Egger
yielded significant results while Brumpton and MR-Twin did not. In general, IVW tended to yield
much stronger p-values than other methods and the family-based methods (Brumpton and MRTwin) tended to be conservative (Supplemental Table S2), in line with our simulation results. In
particular, of the 121 usable trait pairs, IVW identified 78 as significant, Egger identified 56 as
significant, Brumpton identified 20 as significant, and MR-Twin identified 19 as significant.

Discussion

We introduced MR-Twin, a method for testing causal effects between pairs of traits within a 248 Mendelian Randomization (MR) framework, which is provably robust to confounding of any strength resulting from population stratification. Our primary contributions are the following: (i) develop-250 ing a digital twin test, originally proposed by Bates et al. (Bates et al., 2020) in the context of 251 genetic association studies, for MR, coupled with a novel statistic for this test; (ii) demonstrating that, by leveraging trio data, our proposed framework is robust to confounding due to population 253 stratification and to biases from the inclusion of genetic instruments with weak effects; (iii) ex-254 tending our framework to the setting of sibling data, a setting not considered by Bates et al; (iv) 255 conducting the first (to our knowledge) large-scale evaluation of the digital twin test framework in 256 comparison with existing methods for MR. We demonstrated that existing MR methods, including 257 those designed to correct for confounding resulting from horizontal pleiotropy, are prone to false 258 positives when there is confounding from population stratification. 250

While population stratification was the focus of this paper, the MR-Twin framework also pro-260 vides immunity to several other types of confounding effects. Theory dictates that MR-Twin is 261 immune to confounding from familial effects such as assortative mating and dynastic effects since 262 these effects operate through the parental genotypes (see Figure 1 of (Brumpton et al., 2020)), 263 though we do not explicitly test this in this manuscript. As recently demonstrated, cross-trait 264 assortative mating is pervasive and impacts many common genetic analyses (Border et al., 2022) including MR (Hartwig et al., 2018), so this represents another valuable aspect of MR-Twin even 266 if population stratification is believed to be well-controlled in a particular study. In general, MR-267 Twin is immune to any confounder that is independent of the genotypes of offspring given the genotypes of their parents. We note that when we refer to "immunity" we mean this in a theoretical sense – for instance, under the assumption that the model for Mendelian inheritance is correct.

In our particular implementation, we assume that the genetic instruments have been selected to be roughly independent and thus we can sample digital twin genotypes from the parental genotypes using a binomial model. In practice, of course, genetic variants on the same chromosome are never perfectly independent, though with appropriate caution the dependence is weak enough that the effect on calibration should be negligible. More complex models of meiosis will also rely on other factors such as haplotype phasing accuracy.

In addition to population and familial effects, we highlight two under-appreciated sources of bias 277 in MR studies, both of which MR-Twin avoids without requiring the user to modify any parameters 278 or arguments. The first is weak instrument bias (Burgess et al., 2011), which can bias the effect 279 estimate of standard MR methods, including the Brumpton approach (Brumpton et al., 2020). 280 This accounts for the Brumpton method yielding inflated FPRs when the confounding effects were 281 strong (Figure 2A). One of the most common ways to control for weak instrument bias is by filtering 282 out variants with a weak association signal, often with a threshold of F < 10 for the association 283 between a variant and the exposure trait. However, this procedure has been criticized (Burgess 284 et al., 2011) and may not fully correct for weak instrument bias. Other MR methods may also 285 be affected by this bias. In two-sample study designs, the direction of the bias is towards the null 286 rather than the confounded exposure-outcome association estimate (Lawlor, 2016), but the bias 287 remains. 288

Additionally, we found that standard MR methods (IVW, Egger, Median, Mode, etc) may 289 have inflated FPR when there is population structure that induces correlation between variants, 290 even in the absence of stratification (Figure S2). The reason for the induced correlation is that, 291 even though the variants were simulated independently, they were correlated with one another 292 through the population labels. For example, suppose we have two variants, X1 and X2 and a 293 population label U. The causal diagram for these three variables is $X1 \leftarrow U \rightarrow X2$, so X1 and X2 294 are correlated. Our findings corroborate earlier findings that correlation between SNPs can cause 295 calibration issues in MR methods (Burgess et al., 2013). This phenomenon should be taken into account when performing MR simulations or when applying MR to real datasets where variants 297 may be correlated. In the latter case, users should obtain SNP correlations from an appropriately 298

population-matched (Peterson et al., 2019) and sufficiently large (Benner et al., 2017) reference panel.

MR-Twin avoids both of these issues, without requiring the user to specify a SNP correlation 301 matrix or apply various approaches to mitigate weak instrument bias. First, both MR-Twin and 302 Brumpton avoid the correlated-variant issue because they condition on parental genotypes, severing 303 the link between the offspring genotypes and the population structure. Second, MR-Twin would not 304 lose FPR calibration due to weak instrument bias, because this phenomenon has nothing to do with 305 the aspects of the MR-Twin test that guarantee immunity from confounding due to population and 306 familial effects (sampling digital twin genotypes conditioned on parental genotypes). Theoretically, it is possible that the bias in the MR effect estimate used in the MR-Twin statistic (Methods) could 308 lower power, but because the MR effect estimate equally affects both the digital twin statistics and 309 the true offspring statistics, a reduction in power seems unlikely and was not observed empirically 310 (Figure S12). 311

There is extensive literature on family-based methods for avoiding confounding due to popula-312 tion structure in genome-wide association studies or linkage analysis (Weiner et al., 2017; Laird and 313 Lange, 2006; Abecasis et al., 2000; Fulker et al., 1999; Thomson, 1995; Spielman et al., 1993). One 314 prominent example is the transmission disequilibrium test (TDT) (Spielman et al., 1993) and the 315 more-recent polygenic TDT (pTDT) (Weiner et al., 2017). Bates et al (Bates et al., 2020) compare 316 the digital twin test (DTT) to the TDT and show that the DTT is a generalization of the TDT 317 and highlight some of its benefits. Because it is not immediately obvious how to adapt the TDT 318 and pTDT to MR, we do not evaluate their potential use in this context. 319

There are several considerations that come into play when applying the MR-Twin method, which 320 we note here. First, the number of digital twins simulated involves a trade-off between speed and 321 precision (Figure S9). While MR-Twin was slower than competing MR methods (Figure S3), it still 322 ran in a few minutes or less per run on both simulated and real data analyses, justifying the use of 323 a fairly large number of digital twins if possible. Consequently, we recommend 1000 or more digital 324 twins for real data analysis, which should be computationally feasible and precise (and, again. 325 parallelization can make this quite efficient). 100 digital twins is likely sufficient in simulations 326 where there are many replicates and speed is the paramount concern. Second, the populations 327 of the external and family datasets should be similar. This is natural for biobanks like the UK Biobank, but can be more challenging when attempting to combine separate datasets. Third, care should be taken to ensure that the normalization method used and covariates controlled-for are similar in the external and trio datasets in order to avoid potential loss of power.

While the genetic trio offspring used in our UK Biobank analysis were all adults (as all partic-332 ipants in this dataset were aged 40-69 at collection time (Bycroft et al., 2018)), other trio datasets 333 may contain young children. This is a potential issue because some commonly-analyzed traits such 334 as height and weight may not have the same relationship in youths or adolescents as they do in 335 adults, and variants that affect these traits may not yet have realized their full effect in the chil-336 dren yet. Dealing with such time-varying exposures in the context of mendelian randomization is 337 an area of ongoing research (Labrecque and Swanson, 2019), and it is not clear how this would 338 impact MR-Twin results. Even when the offspring of the trios are all adults, it may be difficult 339 to adequately sample certain traits. For example, we were not able to perform MR analysis for complex traits such as heart disease, since none of the offspring in our sample had developed heart 341 disease, largely because all offspring in our sample were aged 40-49. 342

We note a few trends seen across many trait pairs in the real data results, reflecting some 343 practical considerations. First, all standard MR methods identified substantially more trait pairs 344 than did either family-based approach. Given our simulation results showing a large power differ-345 ence in the methods when run with different sample sizes (Figure 2) but similar power when run 346 with the same sample size (Figure 3), along with the fact that the UK Biobank has many more 347 unrelated individuals than trios, we believe that this difference is largely due to the difference in 348 the available sample sizes between unrelated and trio data. The number of trios available as part 349 of public datasets is currently relatively small, limiting the power of family- or trio-based methods, 350 but future increases in the number of available trios will lead to increases in the power of MR-Twin 351 and other family-based methods. Second, some traits pairs had quite different results when the 352 exposure and outcome traits were switched. For example, none of the standard MR methods had 353 significant p-values for DBP \rightarrow Weight, but all were significant for Weight \rightarrow DBP (Supplemental 354 Table S2). This may be due to one causal direction being correct while the other is incorrect, but 355 may also be affected by factors such as differences in the heritability and/or polygenicity between the two traits. 357

Several extensions to the methods presented here are also possible. While we explored contin-

358

uous traits in this paper, further work needs to be done to apply MR-Twin to binary phenotypes 359 such as disease labels. First, a different statistic such as binary cross entropy (rather than our negative squared loss statistic) may be more appropriate. Second, the use of the external effect size 361 estimates in the statistic may have to be modified, depending on the regression method used and 362 the interpretation of the estimates. For example, it would be inappropriate to replace the effect 363 size estimates in our statistic with odds ratios produced by logistic regression. Even for linear 364 data, it is possible that a different statistic than the one we proposed would be more powerful in 365 some situations. Finding most-powerful statistics for a given significance threshold is a direction 366 for future work. Future work could also improve upon the sibling mode of MR-Twin by using 367 population-based priors to infer parental genotypes with a greater level of accuracy, thereby ob-368 taining superior control of false positives. This approach could in principle be developed for and 369 applied to more extended pedigrees. 370

In the Digital Twin Test paper (Bates et al., 2020), Bates et al propose using a Hidden Markov 371 Model (HMM) to simulate digital twins from the parental haplotypes, the latter being generated by 372 phasing the parental genotypes. For the simplicity of avoiding this phasing step and due to the fact 373 that genetic instruments in MR studies are usually selected to be roughly independent (Burgess 374 et al., 2013), we used a simpler method for simulating digital twins using binomial draws from the 375 parental genotypes (Methods). However, the variants used may not be independent even if they 376 appear to be (Burgess et al., 2013), or one may wish to include correlated variants to increase 377 power. Extending MR-Twin to perform the HMM-based digital twin simulation could therefore 378 increase power. 379

Finally, a pre-print from Tudball et al. proposes a randomzation-based approach to MR (Tudball 380 et al., 2022) that, while being conceptually similar, differs from MR-Twin in a few practical aspects. 381 First, Tudball et al. (Tudball et al., 2022) do not discuss the use of external summary statistics to 382 increase power, whereas this is a core part of the MR-Twin approach (as well as in the Digital Twin 383 Test of Bates et al. (Bates et al., 2020)). Second, Tudball et al. develop family-based propensity 384 scores for individual SNPs and suggest aggregating them with Fisher's method or another p-value 385 aggregation method, which is substantially different from our proposed sum-of-squares statistic 386 over all SNPs (Methods). Finally, Tudball et al. employ the HMM-based digital twin simulation 387 model, while (as discussed above) we employ the simpler binomial model. Nevertheless, the broad conceptual similarities between the two methods highlight the promise of randomization-based approaches to make MR findings more robust and the value of continued development to extend these approaches to more complex pedigrees.

$_{92}$ Methods

393 The MR-Twin framework

- We first introduce the standard Mendelian Randomization (MR) model, without any confounding. Suppose that for a collection of N individuals we obtain their genotypes at M SNPs, and a phenotypic measure for an exposure trait and an outcome trait. For a given individual n we denote
 the genotype vector as \mathbf{X}_n , the genotype at some SNP j as \mathbf{X}_{nj} , the exposure trait as E_n , and
 the outcome trait as O_n . Let $(\mathbf{X}_n, E_n, O_n)_{n=1}^N$ denote the collection of these genotypes and traits
 over all N individuals, where (\mathbf{X}_n) is an $(N \times M)$ matrix and (E_n) and (O_n) are $(N \times 1)$ vectors.
 Finally, let \mathbf{X} , E, and O refer to the genotype vector, exposure trait, and outcome trait for a generic
 individual.
- MR uses the genetic "instruments" \mathbf{X} to estimate the effect of an "exposure" trait E on an "outcome" trait O. This estimate is valid regardless of any confounder \mathbf{U} of the association between E and O, assuming that the following conditions hold (Lawlor et al., 2008):
- 1. The genetic instrument **X** is significantly associated with the exposure trait E;
- 2. The genetic instrument \mathbf{X} is independent of any variables (such as those in \mathbf{U}) that confound the relationship between E and the outcome trait O;
- 3. The genetic instrument **X** is not associated with O except due to its association with E.
- The latter two criteria can be captured by the independence statement

$$\mathbf{X} \perp \!\!\!\perp O \mid E$$
 (1)

Assuming these conditions hold, and assuming a linear model for the relationships between the genotypes and phenotypes (a typical assumption in MR analyses), we can test the null hypothesis that there is no direct causal effect of E on O,

$$H_0: \beta_{EO} = 0 \tag{2}$$

such as the ratio estimator $\beta_{EO} = \beta_{XO}/\beta_{XE}$ (when a single instrument is used) or by two stage 414 least squares or inverse-variance weighting (when multiple instruments are used) (Burgess et al., 2013). 416 However, in the case where we have residual population stratification, denoted **Z** (Figure 1A), 417 this independence assumption is violated. This is because, using terminology from Pearl's graphical 418 formalism (Pearl, 1995), $\mathbf{X} \leftarrow \mathbf{Z} \rightarrow O$ is a backdoor path between \mathbf{X} and O, so the two are not 419 marginally independent. Conditioning on E fails to block this backdoor path (i.e. see Figure 1A). 420 Residual population stratification generally cannot be controlled for directly, though approaches 421 such as Principal Components Analysis (PCA) and Linear Mixed Models (LMMs) have been used 422 to reduce its effect (Price et al., 2010). 423 MR-Twin (Figure 1C) is a method that uses family-based genetic data to avoid this confounding. 424 Suppose that, corresponding to each individual's genotypes X, we also observe the genotypes P1 425 and **P2** of their parents (we later relax the trio assumption to allow for parent-child duo or sibling 426 data). Let $\mathbf{A} := (\mathbf{P1}, \mathbf{P2})$. According to the graphical criteria for d-separation developed by Pearl 427 (Pearl, 1995), A d-separates X from Z (Figure 1B): 428

where β_{EO} is not obtained by direct regression but rather via instrumental variables estimators

413

$$\mathbf{X} \perp \!\!\!\perp \mathbf{Z} \mid \mathbf{A}$$
 (3)

This means that, assuming X does not affect some unmeasured variable which in turn affects O (i.e. no horizontal pleiotropy),

$$\mathbf{X} \perp \!\!\!\perp O \mid E, \mathbf{A}$$
 (4)

thereby satisfying the MR conditions regardless of any residual population stratification.

As shown by Bates et al (Bates et al., 2020), the Digital Twin Test framework outlined in Algorithm 3 can be used to perform a hypothesis test conditioned on **A**. The resulting test involves computing the test statistic $t^* = t(\boldsymbol{X}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}}$) (we give the statistic used in this paper in "MR-Twin Test Statistic Incorporating External Weights" below). To perform a test, we construct B random samples $(\tilde{\mathbf{X}}_n)$ where each $\tilde{\mathbf{X}}$ is a random sample given \mathbf{A} with the same distribution as \mathbf{X} given \mathbf{A} (such a sample can be easily constructed using Mendelian inheritance; see "Generating Digital Twins" below). We refer to these samples as "digital twins". For each such sample b, we then compute $t_b = t((\tilde{\mathbf{X}}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}})$, representing a null distribution of genotypes conditioned on the parental genotypes. This, in turn, gives us a p-value for $t^* = \frac{1+\mathbf{1}t_b \geq t^*}{1+B}$, where B is the total number of permutations we perform. The MR-Twin test is therefore a kind of conditional randomization test (Bates et al., 2020; Candes et al., 2018).

Importantly, the proposed algorithm can leverage effect size estimates $(\hat{\beta})$ from any external GWAS datasets (even GWAS datasets where such estimates might be biased due to population stratification) while providing valid tests. The proposed algorithm is robust to any external confounder satisfying Equation 3, such as population stratification, assortative mating, and dynastic effects.

Algorithm 3 Outline of MR-Twin

- 1. Input: Effect sizes for SNPs: $\hat{\beta}$, trio data $\{(\boldsymbol{X}_n, \boldsymbol{A}_n, O_n)_{n=1}^N\}$
- 2. Compute the MR-Twin test statistic $t^* = t((\boldsymbol{X}_n; O_n)_{n=1}^N; \boldsymbol{\hat{\beta}})$
- 3. For b = 1 to B:
 - (a) Sample digital twins $\tilde{\boldsymbol{X}}_n$ given their ancestors \boldsymbol{A}_n .
 - (b) Compute the MR-Twin test statistic $t_b = t((\tilde{\boldsymbol{X}}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}})$
- 4. $p = \frac{1+1\{t_b \ge t^*\}}{1+B}$

Output: p-value: p

456

Next, we detail the MR-Twin test statistic, digital twin generation algorithms, and formal proofs of the exchangeability of digital twins with each other and their real counterparts.

450 Conditional randomization test for mendelian randomization

The MR-Twin test is related to the digital twin test (Bates et al., 2020) and likewise is a kind of conditional randomization test (Candes et al., 2018). Like the digital twin test, MR-Twin leverages the fact that offspring genotypes are conditionally independent of "external" confounders such as population structure given the parental genotypes and uses a conditional randomization test to test the weaker, but equivalent, null hypothesis of no effect conditioned upon the parental genotypes.

Let X be a vector of offspring genotypes, and let A be the genotype vectors of the two parents

of the offspring. **A** may be directly observed, as in trio data, or inferred using parent-child duo or sibling data (see "Generating Digital Twins"). Let **Z** be one or more "external" confounders, defined (Bates et al., 2020) as

$$\mathbf{X} \perp \!\!\!\perp \mathbf{Z} \mid \mathbf{A}$$
 (5)

Thus, population structure is an external confounder, while horizontal pleiotropic traits are not.

We therefore have

$$\mathbf{X}|(\mathbf{Z}, \mathbf{A}) \stackrel{d}{=} \mathbf{X}|\mathbf{A} \tag{6}$$

Assuming that all confounders are external and that \mathbf{X} is significantly associated with E, O is independent of \mathbf{X} given \mathbf{A} under the MR null hypothesis that E has no effect on O. This is because \mathbf{X} would not have any effects on O mediated by E (since E does not affect O under the MR null hypothesis), and all paths not through E are blocked by conditioning on \mathbf{A} as shown in Equation 6. We therefore want to test

$$\mathbf{X} \perp \!\!\!\perp O \mid \mathbf{A} \tag{7}$$

If this holds, then we cannot rule out that either **X** has no effect on *E* or *E* has no effect on *O*.

We test this null hypothesis via a conditional randomization test (Candes et al., 2018).

In testing this null hypothesis, it is helpful to be able to leverage SNP effect sizes estimated from large, external datasets (such as publicly released summary statistics for resources like the UK Biobank (Bycroft et al., 2018)), as this will often yield more statistically significant variants and better effect size estimates than those generated using small genetic family datasets. We therefore note that the following property also holds:

$$\mathbf{X} \perp \!\!\!\perp \hat{\beta} \mid \mathbf{A}$$
 (8)

where we use the shorthand $\hat{\beta}$ to refer to the estimated effect sizes of each SNP on the exposure and outcome traits.

476

We construct "digital twins" $\tilde{\mathbf{X}}$ sampled from the parental genotypes via Mendelian inheritance

477 (see "Generating Digital Twins") such that

$$\tilde{\mathbf{X}}|\mathbf{A} \stackrel{d}{=} \mathbf{X}|\mathbf{A} \tag{9}$$

Given equations 7, 8, and 9, we have the following under the null hypothesis:

$$\mathbf{X}|(\mathbf{A},\hat{\beta},O) \stackrel{d}{=} \mathbf{X}|\mathbf{A} \tag{10}$$

$$\tilde{\mathbf{X}}|(\mathbf{A},\hat{\boldsymbol{\beta}},O) \stackrel{d}{=} \tilde{\mathbf{X}}|\mathbf{A} \tag{11}$$

It follows from equations 9, 10, and 11 that the digital twins are exchangeable under the null hypothesis:

$$\tilde{\mathbf{X}}|(\mathbf{A},\hat{\beta},O) \stackrel{d}{=} \mathbf{X}|(\mathbf{A},\hat{\beta},O)$$
(12)

Therefore, given some statistic $T = t((\boldsymbol{X}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}})$, where N is the number of families,

$$T|(\mathbf{A}, \hat{\beta}, O) \stackrel{d}{=} \tilde{T}|(\mathbf{A}, \hat{\beta}, O)$$
(13)

under the null, where $\tilde{T} = t((\tilde{\boldsymbol{X}}_n; O_n)_{n=1}^N; \boldsymbol{\hat{\boldsymbol{\beta}}})$. We can then use the procedure outline in Algorithm 3 to obtain a p-value for this test statistic (Candes et al., 2018).

484 MR-Twin test statistic incorporating external weights

We construct a test statistic based on a negative sum of squares loss when using X to predict Ovia an MR estimate for the effect of E on O. First, we leverage the effect sizes from the external dataset of the genotype on the exposure trait $\hat{\beta}_{XE}$ to obtain the genetically-predicted exposure trait values:

$$\hat{E}_n = \sum_j \hat{\beta}_{XE,n} \mathbf{X}_{nj} \tag{14}$$

for each individual n and SNP j. We then compute the MR estimate for the effect of the exposure trait on the outcome trait, $\hat{\beta}_{EO}$. This estimate may be a conventional Inverse Variance Weighted

(IVW) estimate (Burgess et al., 2013) or various statistics designed to be robust to pleiotropy such as the Egger-based statistic (Bowden et al., 2015), the weighted median statistic (Bowden et al., 2016), or others. We then predict the outcome trait for each individual n as $\hat{O}_n = \hat{\beta}_{EO}\hat{E}_n$. Finally, we compute the negative squared error of these predictions $-\sum_n (\hat{O}_n - O_n)^2$, summed across all individuals. The full statistic is then

$$t((\boldsymbol{X}_n; O_n)_{n=1}^N; \hat{\boldsymbol{\beta}}) = -\sum_n ((\hat{\beta}_{EO} \sum_j (\hat{\beta}_{XE, n} \mathbf{X}_{nj})) - O_n)^2$$
(15)

96 Generating digital twins

We have assumed that trio data is available thus far for simplicity. However, the MR-Twin framework can also be used when parent-child duo data or sibling data are available. Here we discuss the algorithms used to generate digital twins given trio, parent-child duo, or sibling data.

500 Trio and duo modes

512

513

We assume that the SNPs used in the MR instrument are independent, a common assumption when 501 multi-SNP instruments are used in MR (Burgess et al., 2013). Therefore, we separately sample 502 the genotype of each SNP of the digital twin given the parent and/or offspring genotypes at that 503 SNP. Let (\mathbf{D}_n) be the $(N \times M)$ matrix of digital twin genotypes we will sample, corresponding to 504 the true "offspring" genotypes in (\mathbf{X}_n) . Further, let n index some family and j index some SNP, 505 such that $\mathbf{P1}_{nj}$ (for example) is the genotype for one parent in family n at SNP j. If we have both parents available, sampling \mathbf{D}_{nj} is straightforward. Because the SNPs are considered independent, 507 we do not need to know the parental haplotypes. If a parental genotype $\mathbf{P1}_{nj}$ is 0 or 2, respectively, 508 then a 0 or 1, respectively, is inherited by \mathbf{D}_{nj} . If the parent genotype is 1, then either 0 or 1 is inherited with 50% probability each. \mathbf{D}_{nj} inherits alleles from the two parents independently. This 510 can be summarized as 511

$$\mathbf{D}_{nj} \sim Bern(\mathbf{P1}_{nj}/2) + Bern(\mathbf{P2}_{nj}/2) \tag{16}$$

where Bern stands for the Bernoulli distribution, for each family n and SNP j.

If we only have one parent genotype available, then following Bates et al (Bates et al., 2020),

we fix the offspring's haplotype from the unobserved parent and only simulate a random draw from the observed parent's haplotype. If the observed parent is homozygous, then the allele inherited from that parent is fixed as well, so $\mathbf{D}_{nj} = \mathbf{X}_{nj}$. Otherwise, the allele inherited from this parent will be Bern(0.5). In principle, 0.5 could be replaced with some value based on population allele frequencies. Similarly to the above, the model for the allele from the other parent can be written as $Bern(\mathbf{X}_{nj}/2)$. Thus, if the parent is a heterozygote, we have

$$\mathbf{D}_{nj} \sim Bern(1/2) + Bern(\mathbf{X}_{nj}/2) \tag{17}$$

Sibling mode

In the case where we observe sibling genotypes but not the genotypes of their parents, we assessed two potential approaches. In either case, the observed sibling information is used to infer the probabilities of digital twin genotypes based on the fact that the sibling genotypes give information about the probabilities of various parental genotypes. For instance, a child with a 2 genotype at a SNP guarantees that neither parent has a 0 genotype at that SNP, and makes it more likely that the parents have 2 genotypes than 1 genotypes. Most simply, if one sibling has a 2 genotype at a SNP and the other sibling has a 0, then the parents must both be heterozygotes. In all other cases, approximation is needed.

The first approach is straightforward and involves randomly drawing two haplotypes from the 529 observed sibling haplotypes to generate a digital twin. This shuffling approach gives a rough 530 approximation of the likelihood of digital twin genotypes given the information the observed siblings 531 provide. The second approach, described in the Supplemental Materials, involves using the sibling 532 data to infer a distribution over the possible parents, then performing a weighted random draw of 533 digital siblings based on those parents. In practice, we found that the shuffling approach was faster 534 and yielded lower FPR than the probabilistic approach while achieving similar power, so we used 535 the shuffling approach for the results in this paper. 536

537 Software availability

The code implementing the MR-Twin package can be found at: https://github.com/nlapier2/MR-Twin. Scripts and instructions for repeating the experiments in this paper can be found at: https://github.com/nlapier2/MRTwin-replication. This code is also available in the Supplemental Code files. Please note that UK Biobank genotypes are not publicly released, so those wishing to replicate the experiments will first have to get access to that data via https://www.ukbiobank.ac.uk/enable-your-research/apply-for-access.

Competing interest statement

The authors declare no competing interests.

546 Acknowledgements

- This research has been conducted using the UK Biobank Resource under application number 33127.
- 548 B.F. and S.S. were supported in part by NIH R35GM125055 and NSF CAREER-1943497, III-
- 549 2106908. E.E. and N.L. are funded by NSF award 2106908 and NIH awards U01HG011715 and
- R56HG010812. The authors would also like to thank Matthew J. Tudball for productive discussions
- on potential future work and our respective efforts to develop family-based Mendelian Randomiza-
- 552 tion methods.
- Author contributions: S.S. and E.E. conceived of and supervised the project. N.L., B.F., and
- 554 S.S. developed the methods and wrote the manuscript. N.L., B.F., and S.T. wrote the software
- code and performed the analyses. All authors read and approved the final manuscript.

556 References

- Abecasis GR, Cardon LR, and Cookson W. 2000. A general test of association for quantitative traits in nuclear families. *The American Journal of Human Genetics* **66**: 279–292.
- Balding DJ and Nichols RA. 1995. A method for quantifying differentiation between populations
- at multi-allelic loci and its implications for investigating identity and paternity. Genetica 96:
- 561 3–12.
- Bates S, Sesia M, Sabatti C, and Candès E. 2020. Causal inference in genetic trio studies. Proceed-
- ings of the National Academy of Sciences 117: 24117–24126.

- Benner C, Havulinna AS, Järvelin MR, Salomaa V, Ripatti S, and Pirinen M. 2017. Prospects
 of fine-mapping trait-associated genomic regions by using summary statistics from genome-wide
 association studies. The American Journal of Human Genetics 101: 539–551.
- Berg JJ, Harpak A, Sinnott-Armstrong N, Joergensen AM, Mostafavi H, Field Y, Boyle EA, Zhang
 X, Racimo F, Pritchard JK, et al.. 2019. Reduced signal for polygenic adaptation of height in
 uk biobank. *eLife* 8: e39725.
- Border R, Athanasiadis G, Buil A, Schork AJ, Cai N, Young AI, Werge T, Flint J, Kendler KS,
 Sankararaman S, et al.. 2022. Cross-trait assortative mating is widespread and inflates genetic
 correlation estimates. *Science* 378: 754–761.
- Bowden J, Davey Smith G, and Burgess S. 2015. Mendelian randomization with invalid instruments: effect estimation and bias detection through egger regression. *International Journal of* Epidemiology 44: 512–525.
- Bowden J, Davey Smith G, Haycock PC, and Burgess S. 2016. Consistent estimation in mendelian randomization with some invalid instruments using a weighted median estimator. *Genetic Epi*demiology **40**: 304–314.
- Broadbent JR, Foley CN, Grant AJ, Mason AM, Staley JR, and Burgess S. 2020. Mendelianrandomization v0.5.0: updates to an r package for performing mendelian randomization analyses
 using summarized data. Wellcome Open Research 5.
- Brumpton B, Sanderson E, Heilbron K, Hartwig FP, Harrison S, Vie GÅ, Cho Y, Howe LD, Hughes
 A, Boomsma DI, et al.. 2020. Avoiding dynastic, assortative mating, and population stratification
 biases in mendelian randomization through within-family analyses. *Nature Communications* 11:
 1–13.
- Burgess S, Butterworth A, and Thompson SG. 2013. Mendelian randomization analysis with multiple genetic variants using summarized data. *Genetic Epidemiology* **37**: 658–665.
- Burgess S, Thompson SG, and Collaboration CCG. 2011. Avoiding bias from weak instruments in mendelian randomization studies. *International Journal of Epidemiology* **40**: 755–764.

- Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D, Delaneau
- O, O'Connell J, et al.. 2018. The uk biobank resource with deep phenotyping and genomic data.
- 592 Nature **562**: 203–209.
- ⁵⁹³ Candes E, Fan Y, Janson L, and Lv J. 2018. Panning for gold: 'model-x' knockoffs for high
- dimensional controlled variable selection. Journal of the Royal Statistical Society: Series B
- (Statistical Methodology) **80**: 551–577.
- ⁵⁹⁶ Chen G, Yuan A, Shriner D, Tekola-Ayele F, Zhou J, Bentley AR, Zhou Y, Wang C, Newport MJ,
- Adeyemo A, et al.. 2015. An improved fst estimator. *PLOS One* **10**: e0135368.
- ⁵⁹⁸ Cinelli C, LaPierre N, Hill BL, Sankararaman S, and Eskin E. 2022. Robust mendelian randomiza-
- tion in the presence of residual population stratification, batch effects and horizontal pleiotropy.
- Nature Communications 13: 1093.
- 601 Conomos MP, Reiner AP, Weir BS, and Thornton TA. 2016. Model-free estimation of recent genetic
- relatedness. The American Journal of Human Genetics 98: 127–148.
- 603 Cook JP, Mahajan A, and Morris AP. 2020. Fine-scale population structure in the uk biobank:
- implications for genome-wide association studies. Human Molecular Genetics 29: 2803–2811.
- Fulker D, Cherny S, Sham P, and Hewitt J. 1999. Combined linkage and association sib-pair
- analysis for quantitative traits. The American Journal of Human Genetics 64: 259–267.
- Haase CL, Tybjærg-Hansen A, Ali Qayyum A, Schou J, Nordestgaard BG, and Frikke-Schmidt R.
- 2012. Lcat, hdl cholesterol and ischemic cardiovascular disease: a mendelian randomization study
- of hdl cholesterol in 54,500 individuals. The Journal of Clinical Endocrinology & Metabolism 97:
- 610 E248-E256.
- 611 Hartwig FP, Davey Smith G, and Bowden J. 2017. Robust inference in summary data mendelian
- randomization via the zero modal pleiotropy assumption. International Journal of Epidemiology
- **46**: 1985–1998.
- 614 Hartwig FP, Davies NM, and Davey Smith G. 2018. Bias in mendelian randomization due to
- assortative mating. Genetic Epidemiology 42: 608–620.

- Haworth S, Mitchell R, Corbin L, Wade KH, Dudding T, Budu-Aggrey A, Carslake D, Hemani G,
- Paternoster L, Smith GD, et al.. 2019. Apparent latent structure within the uk biobank sample
- has implications for epidemiological analysis. Nature Communications 10: 1–9.
- Haycock PC, Burgess S, Nounu A, Zheng J, Okoli GN, Bowden J, Wade KH, Timpson NJ, Evans
- DM, Willeit P, et al.. 2017. Association between telomere length and risk of cancer and non-
- neoplastic diseases: a mendelian randomization study. JAMA Oncology 3: 636–651.
- 622 Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D, Laurin C, Burgess S, Bowden
- J, Langdon R, et al.. 2018. The mr-base platform supports systematic causal inference across
- the human phenome. $eLife\ 7: e34408.$
- 625 Hubisz MJ, Falush D, Stephens M, and Pritchard JK. 2009. Inferring weak population structure
- with the assistance of sample group information. Molecular Ecology Resources 9: 1322–1332.
- 627 Labrecque JA and Swanson SA. 2019. Interpretation and potential biases of mendelian randomiza-
- tion estimates with time-varying exposures. American Journal of Epidemiology 188: 231–238.
- 629 Laird NM and Lange C. 2006. Family-based designs in the age of large-scale gene-association
- studies. Nature Reviews Genetics 7: 385–394.
- 631 Lawlor DA. 2016. Commentary: Two-sample mendelian randomization: opportunities and chal-
- lenges. International Journal of Epidemiology 45: 908.
- Lawlor DA, Harbord RM, Sterne JA, Timpson N, and Davey Smith G. 2008. Mendelian random-
- ization: using genes as instruments for making causal inferences in epidemiology. Statistics in
- 635 *Medicine* **27**: 1133–1163.
- Lyall DM, Celis-Morales C, Ward J, Iliodromiti S, Anderson JJ, Gill JM, Smith DJ, Ntuk UE,
- Mackay DF, Holmes MV, et al.. 2017. Association of body mass index with cardiometabolic
- disease in the uk biobank: a mendelian randomization study. JAMA Cardiology 2: 882–889.
- Ochoa A and Storey JD. 2021. Estimating fst and kinship for arbitrary population structures.
- PLOS Genetics 17: e1009241.
- Pearl J. 1995. Causal diagrams for empirical research. *Biometrika* 82: 669–688.

- Peterson RE, Kuchenbaecker K, Walters RK, Chen CY, Popejoy AB, Periyasamy S, Lam M, Iyegbe
- ⁶⁴³ C, Strawbridge RJ, Brick L, et al.. 2019. Genome-wide association studies in ancestrally diverse
- populations: opportunities, methods, pitfalls, and recommendations. Cell 179: 589–603.
- ⁶⁴⁵ Pingault JB, O'reilly PF, Schoeler T, Ploubidis GB, Rijsdijk F, and Dudbridge F. 2018. Using
- genetic data to strengthen causal inference in observational research. Nature Reviews Genetics
- **19**: 566–580.
- eas Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, and Reich D. 2006. Princi-
- pal components analysis corrects for stratification in genome-wide association studies. Nature
- 650 Genetics **38**: 904–909.
- Price AL, Zaitlen NA, Reich D, and Patterson N. 2010. New approaches to population stratification
- in genome-wide association studies. Nature Reviews Genetics 11: 459–463.
- 653 Pritchard JK, Stephens M, and Donnelly P. 2000. Inference of population structure using multilocus
- genotype data. *Genetics* **155**: 945–959.
- Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, De Bakker
- PI, Daly MJ, et al.. 2007. Plink: a tool set for whole-genome association and population-based
- linkage analyses. The American Journal of Human Genetics 81: 559–575.
- ⁶⁵⁸ Spielman RS, McGinnis RE, and Ewens WJ. 1993. Transmission test for linkage disequilibrium:
- the insulin gene region and insulin-dependent diabetes mellitus (iddm). The American Journal
- of Human Genetics **52**: 506–516.
- Thomson G. 1995. Mapping disease genes: family-based association studies. The American Journal
- of Human Genetics **57**: 487.
- 663 Tudball MJ, Smith GD, and Zhao Q. 2022. Almost exact mendelian randomization. arXiv preprint
- arXiv:2208.14035.
- Tyrrell J, Jones SE, Beaumont R, Astley CM, Lovell R, Yaghootkar H, Tuke M, Ruth KS, Freathy
- RM, Hirschhorn JN, et al.. 2016. Height, body mass index, and socioeconomic status: mendelian
- randomisation study in uk biobank. BMJ **352**: i582.

- Wade KH, Carslake D, Sattar N, Davey Smith G, and Timpson NJ. 2018. Bmi and mortality in uk biobank: revised estimates using mendelian randomization. *Obesity* **26**: 1796–1806.
- Weiner DJ, Wigdor EM, Ripke S, Walters RK, Kosmicki JA, Grove J, Samocha KE, Goldstein JI,
- Okbay A, Bybjerg-Grauholm J, et al.. 2017. Polygenic transmission disequilibrium confirms that
- common and rare variation act additively to create risk for autism spectrum disorders. Nature
- 673 Genetics **49**: 978–985.
- Yavorska OO and Burgess S. 2017. Mendelianrandomization: an r package for performing mendelian
- randomization analyses using summarized data. International Journal of Epidemiology 46: 1734—
- 676 1739.

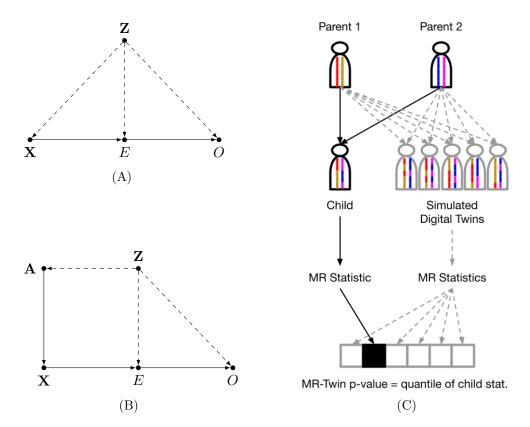


Figure 1: Illustrations of Mendelian Randomization assumptions and the MR-Twin framework. (A) Directed Acyclic Graph (DAG) depicting variables and their relationships in a typical Mendelian Randomization (MR) study, where X is the genotypic instrument, E is the exposure trait, and O is the outcome trait. An external confounder Z, such as population stratification, can cause violations of the MR assumptions. (B) If we have the parental haplotypes A, then X is independent of Z given A. (C) Illustration of the MR-Twin workflow. Digital twin genotypes are sampled from the parental genotypes. MR-Twin is a conditional randomization test, conditioned on A and therefore immune to confounding from Z, in which the p-value is computed based on the quantile of the true offspring's MR-Twin statistic compared to the digital twins' statistics.

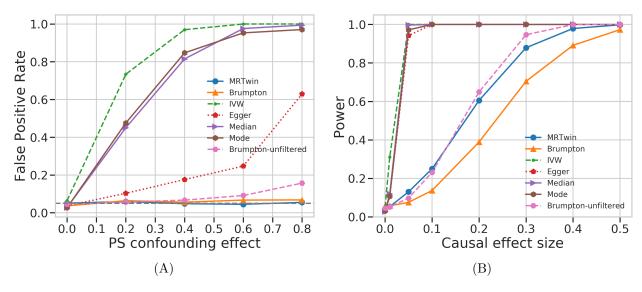


Figure 2: False Positive Rate (FPR) and Power comparison between various methods run on simulated data. (A) False positive rate (y-axis) under varying levels of confounding due to population stratification (PS), with the x-axis describing the magnitude of the confounding effect of population labels on the exposure and outcome trait. (B) Power (y-axis) as a function of the magnitude of the causal effect of the exposure on the outcome trait (x-axis) in a setting with no confounding. Results are averaged over 1000 simulation replicates.

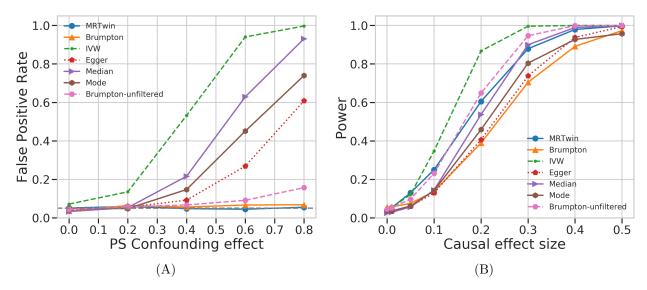


Figure 3: False Positive Rate (FPR) and Power comparison between various methods run on simulated trio data. This is similar to Figure 2 except that IVW, Egger, Median, and Mode are run on the offsprings of the trio dataset instead of the large "external" group of unrelated individuals, such that all methods have the same sample size. (A) False positive rate (y-axis) under varying levels of confounding due to population stratification (PS), with the x-axis describing the magnitude of the effect of the population labels on the exposure and outcome trait. (B) Power (y-axis) as a function of the causal effect size (x-axis). Results are averaged over 1000 simulation replicates.

	MR P-Values			MR-Twin P-Value
Traits	IVW	Egger	Brumpton	MR-Twin
$\overline{\text{LDL Chol.} \rightarrow \text{Total Chol.}}$	$< 10^{-300}$	$< 10^{-300}$	1.64×10^{-11}	$\leq 9.99\times 10^{-4}$
Weight \rightarrow BMI	$< 10^{-300}$	$< 10^{-300}$	4.80×10^{-6}	$\leq 9.99 imes 10^{-4}$
$\mathrm{BMI} \to \mathrm{DBP}$	$\boldsymbol{2.24 \times 10^{-26}}$	5.64×10^{-1}	3.46×10^{-2}	2.69×10^{-1}
$\mathrm{BMI} \to \mathrm{TDI}$	$1.18 imes10^{-19}$	$7.53 imes10^{-3}$	9.99×10^{-2}	8.79×10^{-2}
$\mathrm{Glucose} \to \mathrm{TDI}$	1.54×10^{-1}	2.09×10^{-1}	6.61×10^{-1}	1.91×10^{-1}
$\operatorname{Height} \to \operatorname{Body} \operatorname{Fat}$	9.55×10^{-1}	9.83×10^{-2}	6.73×10^{-1}	5.09×10^{-1}

Table 1: Traditional MR results and MR-Twin results on selected trait pairs from the UK Biobank. Bold numbers are significant at p < 0.05. Note that $9.99 \times 10^{-4} = 1/1001$ is the minimum p-value for MR-Twin in this expriment, as 1000 digital twins were generated. Chol. = Cholesterol; BMI = Body Mass Index; DBP = Diastolic Blood Pressure; TDI = Townsend Deprivation Index.