

1 Title: Negative Control Exposures: Causal effect Identifiability and Use in Probabilistic-Bias and
2 Bayesian Analyses with Unmeasured Confounders

3 Authors: Flanders WD, Waller LA, Zhang Q, Getahun D, Silverberg M, Goodman M.

4 *Corresponding Author:* W. Dana Flanders, wflande@emory.edu; phone: 404-633-7766

5 *Author affiliations:* Department of Epidemiology, Rollins School of Public Health, Emory University,
6 Atlanta, GA 30345 (Flanders WD, Zhang Q, Goodman M); Department of Biostatistics and
7 Bioinformatics, Rollins School of Public Health, Emory University, Atlanta, GA (Flanders WD, Waller
8 LA); Department of Research and Evaluation, Kaiser Permanente Southern California, Pasadena, CA
9 (Getahun D); Division of Research, Kaiser Permanente, Northern California, Oakland, CA (Silverberg M)

10 *Funding information:* This work was supported by Contract AD-12-11-4532 from the Patient Centered
11 Outcome Research Institute and Grant R21HD076387 from the Eunice Kennedy Shriver National
12 Institute of Child Health and Human Development.

13 *Conflict of Interest Statement:* Drs. Flanders and Goodman and Ms. Zhang provide consulting services
14 through Epidemiologic Research & Methods, LLC. This company is owned by Dr. Flanders and provides
15 consulting services to clients. None of the work on the present paper was related to that consulting.

16 **Suggested Running Head:** Negative Controls: Effect Identification & Bias Analyses

17 **Keywords:** bias, confounding, negative controls, negative control exposure, probabilistic bias analysis,
18 Bayesian Analysis, adjustment.

19 **Abbreviations:** RR- risk ratio; CE- causal effect.

20 **Publication history:** This manuscript was submitted to *Epidemiology* (03/17/22), and to medRxiv
21 05.24.22.

1 **Abstract:**

2 Probabilistic bias and Bayesian analyses are important tools for bias correction, particularly if required
3 parameters are nonidentifiable. Negative controls are another tool; they can detect confounding and
4 correct for confounders. Our goals are to present conditions that assure identifiability of certain causal
5 effects and to describe and illustrate a probabilistic bias analysis and related Bayesian analysis that use a
6 negative control exposure.

7 Using potential-outcome models, we characterize assumptions needed for identification of causal effects
8 using a dichotomous, negative control exposure when residual confounding exists. We define bias
9 parameters, characterize their relationships with the negative control and with specified causal effects,
10 and describe the corresponding probabilistic-bias and Bayesian analyses.

11 We exemplify analyses using data on hormone therapy and suicide attempts among transgender people.
12 To address possible confounding by healthcare utilization, we used prior Tdap (tetanus-diphtheria-
13 pertussis) vaccination as a negative control exposure. Hormone therapy was weakly associated with risk
14 (risk ratio (RR) = 0.9). The negative control exposure was associated with risk (RR = 1.7), suggesting
15 confounding. Based on an assumed prior distribution for the bias parameter, the 95% simulation interval
16 for the distribution of confounding-adjusted RR was (0.17, 1.64), with median 0.5; the 95% credibility
17 interval was similar.

18 A dichotomous negative control exposure can be used to identify causal effects when a confounder is
19 unmeasured under strong assumptions. More realistically, assumptions can be relaxed and the negative
20 control exposure may prove helpful for probabilistic bias analyses and Bayesian analyses.

21
22

23 **Introduction:**

24 Residual confounding often threatens valid estimation of causal effects, especially absent randomization
25 of exposure. In a potential outcome framework, confounding implies non-exchangeability, defined below
26 as an association of the exposure with the potential outcomes.

27 Numerous approaches can adjust for, or otherwise account for measured confounders, including
28 restriction to a single level of the confounder, control or adjustment by stratification or modelling in the
29 analysis, difference in difference and regression discontinuity analyses, and use of instrumental
30 variables¹.

31 To detect residual confounding, perhaps due to an unmeasured or mis-measured confounder, one can use
32 a *negative control outcome* or *negative control exposure*². A negative control exposure, our focus, is a

1 variable that does not cause the outcome but is associated with the suspected, but unmeasured confounder
2 (detailed in methods). In an early application, Yerushalmy studied the effects of maternal cigarette
3 smoking on birth weight (reprinted³). To detect confounding, he assessed the association of paternal
4 smoking with his offspring's birthweight; this alternative "exposure" was contemporaneously thought not
5 to affect the outcome of interest – a negative control exposure; he observed an association and interpreted
6 it as suggesting confounding. Later work using cotinine levels suggested this use of paternal smoking as
7 a negative control was valid^{2,4}. As another example, Flanders et al. studied the effects of air pollution on
8 emergency department visits for respiratory diseases. They used pollutant levels the day after the outcome
9 had occurred, which could not cause the outcome, as a negative control exposure to detect residual
10 confounding or other bias^{5,6}. A study of influenza vaccination and deaths from influenza⁷ exemplifies
11 negative control outcomes. A strong association of vaccination status with influenza deaths *before* the
12 influenza season (an alternative "outcome" thought to be unaffected by the exposure of interest – a
13 negative control outcome) suggested bias in the estimated effect of vaccination. Lipsitch et al. discussed
14 and formalized these concepts for detection of residual confounding⁸.

15 Much subsequent work goes beyond bias detection to use negative controls to adjust for residual
16 confounding^{2,9}. Flanders et al. used a negative control exposure to partially correct for residual
17 confounding¹⁰. Their approach, however, involved certain distributional assumptions. Other approaches
18 have involved outcome calibration under a rank preservation assumption¹¹, and use of a linear model for
19 the unmeasured confounder with factor analysis^{12,13}. Miao et al. recently provided conditions that, if met,
20 allow identification of causal effects by using two negative controls, which can act as surrogates for the
21 unmeasured confounder(s)¹⁴⁻¹⁶.

22 Assumptions needed for identifiability can be rather strong if a confounder remains unmeasured. For
23 example, in the categorical case, the approach of Miao et al.¹⁴ requires two negative controls that serve as
24 proxies for the unmeasured confounder (say U). They must have several properties including: each proxy
25 has at least as many categories as U, the proxies are independent conditional on U and certain probability
26 matrices have inverses.

27 Probabilistic bias analyses can address residual biases in the effect estimate that remain after conventional
28 analyses¹⁷. For residual confounding, probabilistic bias analyses use substantive knowledge to help
29 formulate a distribution of bias parameters that characterize unobserved associations (specified in
30 methods), apply that distribution to correct conventional effect estimates, such as risk ratios, and produce
31 a distribution of plausible corrected estimates. Our goals are to describe and illustrate a method that uses a
32 negative control exposure to partially correct for confounding and to formulate probabilistic bias analyses.
33 We extend the approach to a fully Bayesian analysis.

1 **Methods:**

2 *Background, Notation, and Definitions*

3 Our specific objectives are to: present and justify conditions sufficient for using a negative control
4 exposure (N) to identify causal effects of an exposure (E) on an outcome (Y) when a confounder (U) is
5 unmeasured; describe probabilistic bias analyses to address confounding that incorporate information
6 from the negative control; describe a related Bayesian formulation (Appendix 2); and, provide R code to
7 implement these analyses. Here, we measure effects with risk ratios observable in a cohort study. These
8 approaches rely on substantive knowledge to inform the choice of the prior distribution of plausible bias
9 parameters.

10 We assume measured confounders, denoted collectively by X , are categorical, or can be adequately
11 approximated as categorical to control confounding (this imposes little restriction other than regularity
12 conditions). All results are conditional on X , but for simplicity that dependence is suppressed in the
13 notation. For example, M_{enx} denotes the number of people at baseline in the cohort with $E = e, N = n,$
14 $X = x$ for $e, n = 0, 1$ and $x = 1, 2, \dots, |X|$ where $|X|$ is the cardinality of X ; but for simplicity we write
15 M_{en} (conditioning on $X = x$ is implicit). Conditional risk, defined as the probability that the outcome
16 occurs ($Y = 1$) during the follow-up period among those with $E = e, N = n$ at baseline in the cohort, is
17 denoted by $R_{en} = p(Y = 1 | E = e, N = n) = E(Y | E = e, N = n)$. We denote the counterfactual outcome
18 and counterfactual risk among those with $E = e, N = n$, if E were set to e' by $Y(e')$ and $R_{en}(e') =$
19 $E[Y(e') | E = 1, N = n]$ for $e, n, e' = 0$ or 1 (Table 1), respectively; $R_{en}(1)$ or $R_{en}(0)$ must be
20 counterfactual. Similarly, $R(e)$ is the counterfactual risk in the population if E were set to e for all.

21 The assumed causal relationships are summarized in Figure 1, a Single World Intervention Template
22 (SWIT)¹⁸. The causal relationships in the SWIT, assumed correct, imply and are consistent with
23 assumptions A1–A4:

24 A1) $N \perp\!\!\!\perp E, Y(e) | U, X$; $Y(n, e) = Y(e)$, conditional independence between N and E, Y ; N has no effect;

25 A2) $E \perp\!\!\!\perp Y(e) | U, X$; conditional exchangeability;

26 A3) Conditional on E, X , we expect the negative control to be associated with Y and unmeasured
27 confounders U ;

28 A4) $R_{en}(e') = R_{en}$ and $R_{en,u}(e') = R_{en,u}$ if $e' = e$; counterfactual-model consistency.

29 Ideal or U-comparable negative control exposures described by Lipsitch et al.⁸ should satisfy
30 assumptions A1 – A4 (Web Appendix S4). However, additional variables that are not ideal or U-
31 comparable negative controls can satisfy A1 – A4, serve as indicators of residual confounding or other

1 bias¹⁰ and be used for the probabilistic bias or Bayesian analyses described here (e.g., Supplemental
2 Figure S3).

3 We consider the causal effects of exposure among those with $E = e, N = n$, denoted by CE_{en} for $e, n =$
4 $0,1$ (Table 1). Using risk ratios and counterfactuals, we express CE_{en} as:

5 1) $CE_{en} = R_{en}(1)/R_{en}(0)$ for $e, n = 0,1$;

6 and, the population average causal effect as:

7 2) $PACE = R(1)/R(0)$.

8 In the remainder of methods we first consider CE_{10} in detail, providing and justifying assumptions under
9 which CE_{10} can be identified. We then introduce a bias parameter to relax the assumption needed for
10 identification and use this parameter as the basis for probabilistic bias analyses. Finally, we consider
11 other causal effects. In Appendix 2 provide a fully Bayesian formulation of our approach.

12 *Identifiability Conditions and Probabilistic Bias Analysis for CE_{10}*

13 We show that CE_{10} is identifiable in a cohort study under an easily-specified, but strong assumption
14 involving the distribution of the negative control N .

15 Consistent with the pattern of causal effects summarized in Figure 1 and assumptions (A1)-(A4), we can
16 write the identifiable risk R_{00} as:

17 3) $R_{00} = \frac{\sum_u R_{00u} P(E=0|U=u) P(N=0|U=u) p(U=u)}{\sum_u p(E=0|U=u) p(N=0|U=u) p(U=u)}$

18 where R_{enu} is the conditional risk among those with $E = e, N = n, U = u$. Let $R_{enu}(e')$ denote the
19 counterfactual risk among those with $E = e, N = n, U = u$ and E was set to e' . We can write the
20 counterfactual risk $R_{10}(0)$ as the weighted average:

21 4) $R_{10}(0) = \frac{\sum_u R_{10u}(0) P(E=1|U=u) P(N=0|U=u) p(U=u)}{\sum_u p(E=1|U=u) p(N=0|U=u) p(U=u)}$
22 $= \frac{\sum_u R_{00u}(0) P(E=1|U=u) P(N=0|U=u) p(U=u)}{\sum_u p(E=1|U=u) p(N=0|U=u) p(U=u)}$ (substitute $R_{00u}(0)$ for $R_{10u}(0)$, assumption A2)
23 $= \frac{\sum_u R_{01u}(0) P(E=1|U=u) P(N=0|U=u) p(U=u)}{\sum_u p(E=1|U=u) p(N=0|U=u) p(U=u)}$ (substitute $R_{01u}(0)$ for $R_{00u}(0)$ by A1)
24 $= \frac{\sum_u R_{01u} P(E=1|U=u) P(N=0|U=u) p(U=u)}{\sum_u p(E=1|U=u) p(N=0|U=u) p(U=u)}$ (substitute R_{01u} for $R_{01u}(0)$ by consistency)

25 We now state two assumptions either of which, with Assumptions (A1-A4), suffices (Claim 1) to assure
26 identifiability of CE_{10} :

1 A5a) $p(E = 1|U = u) = p(N = 1|U = u)$ for all u ; (equality of conditional distributions) or

2 A5b) $p(U = u|E = 1, N = 0) = p(U = u|E = 0, N = 1)$, for all u .

3 Claim 1: Under assumptions (A1-A4) and (A5a) or (A5b): (i) $R_{10}(0) = R_{01}$; and (ii) CE_{10} is identified by
4 the ratio of observable risks R_{10}/R_{01} .

5 Proof: By assumption (A5a), we can substitute $p(N = 1|U = u)$ for $p(E = 1|U = u)$, and $p(E = 0|U =$
6 $u)$ for $p(N = 0|U = u)$ into the last line of Expression (4), showing that $R_{10}(0) =$

7 $\frac{\sum_u R_{01u} P(E=0|U=u) P(N=1|U=u) p(U=u)}{\sum_u p(E=0|U=u) p(N=1|U=u) p(U=u)}$ which equals R_{01} proving (i). Proof of Claim (i) using (A5b) is

8 similar (Web Appendix S5). Now, $CE_{10} = \frac{R_{10}}{R_{10}(0)}$ which, by (i), equals $\frac{R_{10}}{R_{01}}$. The latter is consistently

9 estimated by the ratio of observable risks in the appropriate subgroups of the cohort study, proving (ii).

10 Following Claim 1, we take the identifiable risk ratio R_{10}/R_{01} as the estimator of CE_{10} .

11 Note: The intuition behind this estimator is that the distortion caused by the association of the
12 unmeasured variable U with exposure – is compensated for and balanced by the association of U with the
13 negative control; under assumption (A5b), the distribution of U is the same in the groups being compared.

14 Assumptions (A5a-A5b) differ from the equi-distributional confounding assumption of Sofer et al.¹⁹
15 which concerns equality of the conditional distributions of the outcome and of a negative control outcome
16 (NCO); see also the “confounding bridge” assumption of Miao et al.²⁰ that involves conditional
17 distributions of the NCO and the outcome, rather than the negative control exposure and the exposure.

18 There is some plausibility that assumption (A5a) or (A5b) would hold, at least approximately, since
19 negative controls “... should be selected such that they share a common confounding mechanism as the
20 exposure and outcome variables ...”². Nevertheless, the assumption (A5a) or (A5b) is strong and, with U
21 unmeasured, unverifiable. Therefore, we introduce a bias parameter that allows for deviations from (A5a-
22 A5b) and that can be used in probabilistic bias analyses. In particular, we relax the key implication of
23 assumption (A5a, A5b) that $R_{10}(0) = R_{01}$, and instead assume:

24 A5c) $\frac{R_{10}(0)}{1-R_{10}(0)} = \varepsilon_1 \frac{R_{01}}{1-R_{01}}$.

25 A5c assumes that the counterfactual odds $R_{10}(0)$ is equal to the (observable) odds $R_{01}/(1 - R_{01})$ times a
26 bias parameter ε_1 . Using risk odds in assumption (A5c), rather than say risks, assures that $R_{10}(0)$ is not
27 outside the range $(0,1)$ for $\varepsilon_1 \in (0, \infty)$. Substituting $\varepsilon_1 R_{01}/(1 - R_{01} + \varepsilon_1 R_{01})$ for $R_{10}(0)$, justified by
28 Assumption (A5c), gives:

$$1 \quad 5) \quad CE_{10} = \frac{R_{10}}{R_{10}(0)} = \frac{1}{\varepsilon_1} \frac{R_{10}(1-R_{01}+\varepsilon_1 R_{01})}{R_{01}}$$

2 Note: The bias parameter ε_1 reflects residual bias in R_{10} as an estimator of the counterfactual risk $R_{10}(0)$
3 on the odds scale. Equation (5) shows that ε_1 equals the ratio of the estimator to the causal effect,
4 sometimes referred to as the confounding risk ratio^{1,21-23}, multiplied by $(1 - R_{01}(1 - \varepsilon_1))$. For rare
5 outcomes, $(1 - R_{01}(1 - \varepsilon_1)) \approx 1$, and ε_1 approximates the confounding risk ratio. This
6 conceptualization may aid interpretation (see also, Supplemental Web Appendix S1).

7 To implement a probabilistic bias analysis for residual confounding, we specify a distribution for ε_1 (with
8 support from 0 to infinity); the distribution has greater or less weight in the tails, depending on the extent
9 to which $R_{10}(0)$ and R_{01} are thought to differ (reflecting differences in the conditional distributions of E
10 and N). If the negative-control association with U is thought to mirror the corresponding exposure
11 association fairly accurately, the ε_1 -distribution can be formulated with a substantial probability that ε_1 is
12 near 1, whereas if the associations differ substantially, greater weight can be assigned elsewhere. An R
13 program to implement probabilistic bias analysis is given in Web Appendix S2. Using a Monte Carlo
14 approach, the program: randomly selects a value of the bias parameter from the specified distribution;
15 applies the bias model to calculate the counterfactual risk ($R_{10}(0)$, Assumption A5c); accounts for
16 random error by sampling $R_{10}(0)$ and R_{10} from binomial distributions; applies Equation (5) to calculate a
17 bias-adjusted estimate of CE_{10} ; and creates a simulation interval.¹⁷

18 *Identifiability Conditions and Probabilistic Bias Analysis for CE_{01}*

19 Assumptions (A1-A4; A5a or A5b) imply that $CE_{01} = CE_{10}$ (proven as Claim 2 in Web Appendix S5),
20 which implies that CE_{01} , like CE_{10} , is identifiable as the ratio of observable risks R_{10}/R_{01} . To relax
21 assumption (A5a or A5b) and conduct probabilistic bias analyses, we introduce a second bias parameter
22 (ε_2), that plays a role like that of ε_1 : ε_2 reflects differences between the estimator $\frac{R_{10}}{R_{01}}$ based on observed
23 risks and the estimand CE_{01} . We have assumption (A5d):

$$24 \quad A5d) \quad \frac{R_{01}(1)}{1-R_{01}(1)} = \varepsilon_2 \frac{R_{10}}{1-R_{10}}$$

25 implying:

$$26 \quad 6) \quad CE_{01} = \frac{R_{01}(1)}{R_{01}} = \frac{R_{10}}{R_{01}(\varepsilon_2 - \varepsilon_2 R_{10} + R_{10})}$$

27 By specifying a distribution for ε_2 we can conduct probabilistic bias analyses for CE_{01} , like those for
28 CE_{10} .

29 *Identifiability Conditions and Probabilistic Bias Analysis for CE_{11} , CE_{00} and PACE*

1 Causal effects CE_{11} , CE_{00} and $PACE$ can differ from CE_{10} and CE_{01} and so are not necessarily identified
2 as R_{10}/R_{01} under Assumptions A1-A4, A5b. However, we can identify these effects if we can assume a
3 multiplicative model for the effect of E given U , conditional on $U = u$ and $N = n$:

4 A6) $R_{enu}(e') = e^{\beta_1 \cdot e'} R_{0nu}$ (multiplicative homogeneity of effects E)

5 In words, assumption (A6) states that the counterfactual risk, if E were set to e' , is $e^{\beta_1 \cdot e'}$ times the risk
6 among those with $E = 0$, $N = n$ and $U = u$. We show in Appendix 1 (Claim 3) that Assumption (A6)
7 implies that CE_{en} and $PACE$ both equal $\frac{R_{en(1)}}{R_{en(0)}} = e^{\beta_1}$. By Assumptions (A1-A5), CE_{10} (and CE_{01}) are
8 identified, and therefore so are CE_{11} , CE_{00} and $PACE$. Using bias parameters to account for errors in
9 assumptions (A1-A6) and combining results, we have

10 A7) $CE_{11} = \frac{1}{\varepsilon_3} CE_{10}$ and $CE_{00} = \frac{1}{\varepsilon_4} CE_{10}$.

11 Analyses based on use of ε_3 and ε_4 use the strong assumption of a multiplicative effect of exposure in
12 addition to (A1-A5) and are viewed as supplementary.

13 *Fully Bayesian Analysis for CE_{10} and CE_{01}*

14 Appendix 2 outlines use of prior distributions for R_{10} and R_{01} and two parameters κ_1 and κ_2 to provide a
15 fully Bayesian formulation of the problem²⁴. Parameters κ_1 and κ_2 (defined in Appendix 2: Equation 1A
16 and just below) reflect the same associations in the Bayesian formulation as do ε_1 and ε_2 in probabilistic
17 bias analyses. Web Appendix S3 documents an R program to implement Bayesian analyses for CE_{10} and
18 CE_{01} .

19 **Example**

20 We illustrate these methods by applying them to investigate the possible effect of gender-affirming
21 hormone therapy on risk of suicide attempts. We use data from ongoing studies of transgender people
22 ^{25,26}, summarized in Table 2. The cohort consists of people from two health plans in California, who were
23 20 years old or younger on December 31, 2015, and received a transgender-specific diagnosis (e.g.,
24 'gender dysphoria') by age 20. We defined exposure as receiving gender-affirming hormones or puberty
25 suppression therapy at or before age 20. The outcome of interest was at least one episode of self-inflicted
26 injury or poisoning, or any hospitalization or emergency room visit for a mental health problem
27 documented in the medical records during the 1-year follow-up period starting at age 20. Web Appendix
28 S6 includes additional descriptive information about the cohort. We were concerned about potential
29 confounding by healthcare utilization, as greater utilization might associate with both more hormone
30 therapy and more (documentation of) mental health diagnoses. Therefore, we used recorded receipt of

1 Tdap vaccine at or before age 20 years as a negative control exposure. If healthcare utilization was a
2 confounder, we thought that Tdap vaccination should, like hormone therapy, be associated with both
3 healthcare utilization and the outcome. The crude risk ratio (cRR) for exposure is 0.88. After adjusting for
4 the negative control, the Mantel-Haenszel mRR is 0.88 (95% CI: 0.56 – 1.38; Table 3). The negative
5 control was associated with risk, both among the exposed ($RR_1 = 1.70$; 95% CI: 0.76 – 3.79) and the
6 unexposed ($RR_2 = 1.66$; 95% CI: 1.07 – 2.56).

7 To implement the probabilistic bias analyses, we specify prior distributions for ε_1 . As provided in the R
8 program (Web Appendix S2), we chose a log-normal distribution for ε_1 , with median 1 and the ratio of
9 the 10th percentile to the median of 0.5. For this specification, the 10th and 90th percentiles of the prior for
10 ε_1 are 0.50 and 2.0, indicating that ε_1 will fall in this range with 80% probability under the prior. The
11 resulting distribution of confounding-adjusted causal effect estimates (simulation intervals)¹⁷ is given in
12 Figure 2, for the assumed prior. The 95% simulation interval is from 0.17 to 1.64, with median 0.52.
13 These results are approximately interpretable as semi-Bayesian, as the bias parameter is sampled from a
14 prior distribution¹⁷. Using the fully Bayesian analysis described in Appendix 2 with uninformative priors
15 for $P(Y = 1|E = 1, N = 2)$ and $P(Y = 1|E = 2, N = 1)$, a log-normal prior for κ_1 and setting the
16 median of the prior (log-normal) to 1 and the variance so that the ratio of the median to the 10th percentile
17 was 0.5, the 95% credibility interval was (0.19 to 1.69) with median 0.54. In a sensitivity analysis, we
18 doubled the variance of the prior for κ_1 - reflecting greater uncertainty in the value of the bias parameter.
19 The 95% credibility interval was then (0.16, 2.08).

20

21 Discussion

22 We have justified assumptions that, if correct, imply one can use a negative control exposure to account
23 for unmeasured confounding and identify causal effects (CE_{10} and CE_{01}). If exposure effects are
24 multiplicative, this suffices to also identify CE_{11} and CE_{00} as well as other effects, such as the population
25 average effect. The key assumption (A5a or A5b) is strong, so we have also described and illustrated two
26 ways to relax the assumption. These inter-related methods, probabilistic bias analyses and Bayesian
27 analyses, both use information from a negative control exposure to account for residual confounding and
28 require the researcher to specify a prior distribution for a bias parameter; they yielded similar results, as
29 expected, in our example. Our formulation of probabilistic bias analyses, as is common¹⁷, includes a bias
30 model, postulating a prior distribution for the bias parameters, and Monte Carlo simulation to obtain a
31 distribution for the bias-adjusted estimate of interest. Our method extends the approach to incorporate
32 information from the negative control. We also describe a complementary Bayesian formulation.

1 Exchangeability implies that risk among those actually exposed is the same as the risk among the
2 unexposed if they had been exposed, and conversely; it can be defined as independence of the actual
3 exposure and response types²⁷⁻²⁹. The methods proposed for addressing non-exchangeability are natural
4 ones in the sense that they use the negative control as a reflection of exposure associations with
5 unmeasured confounders that define non-exchangeability. If the negative control has more than two
6 categories, say $N = n$ for $n = 0, 1, \dots, N$, then the approach described here is still applicable by selecting
7 two categories (or combinations of categories) and contrasting them. For example, if a priori knowledge
8 suggested that $P(U = u|E = 1, N = 0) \approx P(U = u|E = 0, N \in S_1)$ where $S_1 = \{3, 4\}$, then $R_{10}(0)$ could
9 be estimated by R_{0S_1} . To the extent allowed by *a priori* knowledge, original categories can be combined
10 to most closely approximate the identifying assumption.

11 Two causal effects are most readily addressed by the proposed approach, the effect of exposure among
12 the exposed without negative control exposure ($E = 1, N = 0$), and the effect of exposure among the
13 unexposed with negative control exposure ($E = 0, N = 1$). The researcher can use substantive knowledge
14 to assess how well the association of the negative control with the unmeasured confounder reflects the
15 association of exposure with the confounder. In the example when using Bayesian analyses, doubling the
16 variance of the prior distribution for κ_1 led to wider credibility intervals, but not substantially so
17 (Example 1) – suggesting some degree of robustness to a modest change in uncertainty regarding the prior
18 distribution of κ_1 . While it is possible to extend the approach to apply it to effects in other subgroups and
19 to a population average causal effect (Appendix 1), assignment of the prior distribution is perhaps more
20 uncertain because an additional assumption (multiplicative effect of exposure) is needed. Therefore, we
21 view these additional analyses as secondary. We caution that an association between a negative control
22 and the outcome can reflect bias other than non-exchangeability, such model mis-specification⁵. Our
23 analyses are not designed to correct for these other biases.

24 Results of the probabilistic bias analysis and the Bayesian analysis both utilize a negative control to
25 correct for residual confounding and are based on researcher-supplied inputs that are informed, to the
26 extent possible, on subject matter knowledge. The probabilistic bias analyses depend on the prior
27 distributions for ε_1 and ε_2 and the Bayesian results on those for κ_1 and κ_2 ; here we used log-normal
28 distributions for both. With a non-informative prior for the other parameters, the 95% credibility interval
29 (Bayesian analysis) was similar to the simulation interval (probabilistic bias analysis). Absent
30 assumptions (such as A5b), parameters κ_1 and κ_2 or bias parameters ε_1 and ε_2 are not identifiable;
31 however, “indirect learning” is possible³⁰, evidenced here in the change from prior to posterior
32 distribution of κ_1 (Web Appendix S1: Figures S1 and S2). This indirect learning and changes in the
33 distributions result from learning about the identifiable parameters³⁰.

1 In some situations, when the confounder is known but unmeasured, external results may provide direct
2 estimates of κ_1 or ε_1 . For example using equation 5, ε_1 could be estimated as $\varepsilon_1 = \frac{R_{10} - R_{10}R_{01}}{aCE_{10}R_{01} - R_{10}R_{01}}$
3 where aCE_{10} is an external estimate of the causal effect (e.g., adjusted for all confounders in a study
4 where U was measured). However, if external measurements of the exposure, confounders and outcome
5 are available, we can also consider other, possibly more efficient, approaches¹⁷ or perhaps use of a
6 directly calculated confounding risk ratio. We could also use priors for both κ_1 and a causal effect (e.g.,
7 CE_{01}), plus another parameter (e.g., R_{10}). Evaluation of the posterior distribution, however, would likely
8 then require Gibbs sampling or other technique more complicated than the straightforward one used here.

9 In summary, we have provided assumptions sufficient for using a negative control to identify causal
10 effects when a confounder is unmeasured, and have described and illustrated the application of both
11 probabilistic bias analysis and Bayesian formulations to address residual confounding. The latter methods
12 use a negative control exposure, use researcher-supplied prior information about how well the negative
13 control captures the associations that create confounding, and produce results partially adjusted for the
14 residual confounding.

15

1 References

- 2 1. Lash TLV, Tyler J., Haneuse S, Rothman KJ. *Modern epidemiology*. 4th ed. Philadelphia: Wolters
3 Kluwer Health/Lippincott Williams & Wilkins; 2021.
- 4 2. Shi X, Miao W, Tchetgen Tchetgen EJ. A selective review of negative control methods in
5 epidemiology. *Current Epidemiology Reports*. 2020;1-13.
- 6 3. Yerushalmy J. The relationship of parents' cigarette smoking to outcome of pregnancy—
7 implications as to the problem of inferring causation from observed associations. *International*
8 *journal of epidemiology*. 2014;43(5):1355-1366.
- 9 4. Taylor AE, Smith GD, Bares CB, Edwards AC, Munafò MR. Partner smoking and maternal cotinine
10 during pregnancy: implications for negative control methods. *Drug and alcohol dependence*.
11 2014;139:159-163.
- 12 5. Flanders WD, Klein M, Darrow LA, et al. A Method to Detect Residual Confounding in Spatial and
13 Other Observational Studies. *Epidemiology (Cambridge, Mass)*. 2011;22(6):823.
- 14 6. Flanders WD, Klein M, Strickland M, et al. A method of identifying residual confounding and
15 other violations of model assumptions. *Epidemiology*. 2009;20(6):S44-S45.
- 16 7. Jackson LA, Jackson ML, Nelson JC, Neuzil KM, Weiss NS. Evidence of bias in estimates of
17 influenza vaccine effectiveness in seniors. *International journal of epidemiology*. 2006;35(2):337-
18 344.
- 19 8. Lipsitch M, Tchetgen Tchetgen EJ, Cohen T. Negative Controls: A Tool for Detecting Confounding
20 and Bias in Observational Studies. *Epidemiology*. 2010;21(3):383-388.
- 21 9. Swanson SA, Hernán MA, Miller M, Robins JM, Richardson TS. Partial identification of the
22 average treatment effect using instrumental variables: review of methods for binary
23 instruments, treatments, and outcomes. *Journal of the American Statistical Association*.
24 2018;113(522):933-947.
- 25 10. Flanders WD, Strickland MJ, Klein M. A new method for partial correction of residual
26 confounding in time-series and other observational studies. *American journal of epidemiology*.
27 2017;185(10):941-949.
- 28 11. Tchetgen Tchetgen EJ. The control outcome calibration approach for causal inference with
29 unobserved confounding. *American journal of epidemiology*. 2014;179(5):633-640.
- 30 12. Wang J, Zhao Q, Hastie T, Owen AB. Confounder adjustment in multiple hypothesis testing.
31 *Annals of statistics*. 2017;45(5):1863.
- 32 13. Jacob L, Gagnon-Bartsch JA, Speed TP. Correcting gene expression data when neither the
33 unwanted variation nor the factor of interest are observed. *Biostatistics*. 2016;17(1):16-28.
- 34 14. Miao W, Geng Z, Tchetgen Tchetgen EJ. Identifying causal effects with proxy variables of an
35 unmeasured confounder. *Biometrika*. 2018;105(4):987-993.
- 36 15. Shi X, Miao W, Nelson JC, Tchetgen Tchetgen EJ. Multiply robust causal inference with double-
37 negative control adjustment for categorical unmeasured confounding. *Journal of the Royal*
38 *Statistical Society: Series B (Statistical Methodology)*. 2020;82(2):521-540.
- 39 16. Kuroki M, Pearl J. Measurement bias and effect restoration in causal inference. *Biometrika*.
40 2014;101(2):423-437.
- 41 17. Lash TL, VanderWeele TJ, Haneuse S, Rothman KJ. Bias Analysis (Chapter 21). In: *Modern*
42 *Epidemiology*. 4th ed. Philadelphia: Wolters Kluwer Health/Lippincott Williams & Wilkins;
43 2021:711-754.
- 44 18. Richardson TS, Robins JM. Single World Intervention Graphs (SWIGs): A Unication of the
45 Counterfactual and Graphical Approaches to Causality. *Working Paper Number 128, Center for*
46 *Statistics and the Social Sciences, University of Washington*.
47 2013:<http://www.csss.washington.edu/Papers/wp128.pdf>.

- 1 19. Sofer T, Richardson DB, Colicino E, Schwartz J, Tchetgen Tchetgen E. On negative outcome
2 control of unobserved confounding as a generalization of difference-in-differences. *Statistical*
3 *science: a review journal of the Institute of Mathematical Statistics*. 2016;31(3):348.
- 4 20. Miao W, Shi X, Tchetgen Tchetgen EJ. A confounding bridge approach for double negative
5 control inference on causal effects. *arXiv preprint arXiv:180804945*. 2018.
- 6 21. Flanders WD, Khoury M. Indirect assessment of confounding: graphic description and limits on
7 effect of adjusting for covariates. *Epidemiol*. 1990;1(3):239-246.
- 8 22. Yanagawa T. Case-control studies: assessing the effect of a confounding factor. *Biometrika*.
9 1984;71(1):191-194.
- 10 23. Miettinen OS. Components of the crude risk ratio. *American Journal of Epidemiology*.
11 1972;96(2):168-172.
- 12 24. Greenland S. Multiple-bias modelling for analysis of observational data. *Journal of the Royal*
13 *Statistical Society: Series A (Statistics in Society)*. 2005;168(2):267-306.
- 14 25. Mak J, Shires DA, Zhang Q, et al. Suicide attempts among a cohort of transgender and gender
15 diverse people. *American journal of preventive medicine*. 2020;59(4):570-577.
- 16 26. Goodman M, Nash R. Examining health outcomes for people who are transgender. *Washington,*
17 *DC: Patient-Centered Outcomes Research Institute (PCORI)* <https://doi.org/10.25302/22019AD>.
18 2018;12114532.
- 19 27. Flanders WD, Eldridge RC. Summary of relationships between exchangeability, biasing paths and
20 bias. *European journal of epidemiology*. 2015;30(10):1089-1099.
- 21 28. Hernán MA. A definition of causal effect for epidemiology. *J Epidemiol Community Health*.
22 2004;58:265-271.
- 23 29. Greenland S, Robins J. Identifiability, exchangeability, and epidemiologic confounding. *Int J*
24 *Epidemiol*. 1986;15:413-419.
- 25 30. Gustafson P. On model expansion, model contraction, identifiability and prior information: two
26 illustrative scenarios involving mismeasured variables. *Statistical science*. 2005;20(2):111-140.

27

28

1 **Appendix 1**

2 Claim 3: Under assumptions (A1-A4, A5b, and A6), $CE_{en} = R_{en}(1)/R_{en}(0) = e^{\beta_1}$ and $PACE = e^{\beta_1}$.

3 Proof: $R_{en}(e')$ is a weighted average of the counterfactual outcomes $R_{enu}(e')$. Therefore:

$$\begin{aligned}
 4 \quad \frac{R_{en}(1)}{R_{en}(0)} &= \frac{\sum_u R_{enu}(1)P(E=e|U=u)P(N=n|U=u)p(U=u)}{\sum_u p(E=e|U=u)p(N=n|U=u)p(U=u)} / \frac{\sum_u R_{enu}(0)P(E=e|U=u)P(N=n|U=u)p(U=u)}{\sum_u p(E=e|U=u)p(N=n|U=u)p(U=u)} \\
 5 &= \frac{\sum_u R_{enu}(1)P(E=e|U=u)P(N=n|U=u)p(U=u)}{\sum_u R_{enu}(0)P(E=e|U=u)P(N=n|U=u)p(U=u)} \\
 6 &= \frac{\sum_u R_{enu}(0)e^{\beta_1}P(E=e|U=u)P(N=n|U=u)p(U=u)}{\sum_u R_{enu}(0)P(E=e|U=u)P(N=n|U=u)p(U=u)} \\
 7 &= e^{\beta_1}.
 \end{aligned}$$

8 Thus: $CE_{11} = CE_{00} = CE_{10} = CE_{01}$. The $PACE$ is a weighted average of the four E, N – specific effects
 9 CE_{en} (with weights $P(E = e, N = n)$), so $PACE$ equals the common value e^{β_1} .

10

11 **Appendix 2**

12 The Appendix describes a Bayesian formulation of our approach. The analytic goal now is to calculate
 13 Bayesian credibility intervals and the posterior median for the causal effects, conditional on the
 14 observations and using the negative control exposure. We use the terminology and definitions from
 15 Lash^{1,17}. We present details only for the analysis of CE_{10} , as those for CE_{01} are directly analogous.
 16 Results for CE_{11} and CE_{00} depend on an additional, strong assumption and are considered supplementary.

17 Parameters

18 The bias parameter ε_1 of the main text was introduced to relax assumption (A5a); it allows for differences
 19 between the distribution of E given U and that of N given U . Following assumption (A5c), we define the
 20 (relative) bias parameter \varkappa_1 as:

$$21 \quad 1A) \quad \varkappa_1 = \frac{R_{10}(0)}{1-R_{10}(0)} / \frac{R_{01}}{1-R_{01}}.$$

22 Parameters R_{10}, R_{01} (Table 1) and \varkappa_1 fully parameterize the conditional distributions of: outcome Y
 23 and $Y(0)$ among those with $E = 1, N = 0$, and those among individuals with $E = 0, N = 1$. For
 24 example, the distribution of $Y(0)$ among those with $E = 1, N = 0$ is:

$$25 \quad P(Y(0) = 1|E = 1, N = 0) = R_{10}(0) = \varkappa_1 R_{01} / (1 - R_{01} + \varkappa_1 R_{01}), \text{ (expression 1A).}$$

26 CE_{10} can be expressed using only R_{10}, R_{20} and \varkappa_1 . Much like ε_1 , we define $\varkappa_2 = \frac{R_{01}(1)}{1-R_{01}(0)} / \frac{R_{10}}{1-R_{10}}$ so that
 27 $P(Y(1) = 1|E = 0, N = 1) = R_{01}(1) = R_{10}/R_{01}(\varepsilon_2 - \varepsilon_2 R_{10} + R_{10})$. This formulation parallels that of
 28 probabilistic bias analysis in the main text; e.g., Equation 1A replaces Assumption A5a).

1 Sampling Distribution

2 The conditional likelihood of the observed data, given the parameters $R_{10} = p_1$ and $R_{01} = p_2$ is:

3 2A) $P(Y = y_1 | E = 1, N = 0, R_{10} = p_1)P(Y = y_0 | E = 0, N = 1, R_{01} = p_0)$
 4
$$= \frac{M_{10}!}{(M_{10}-y_1)!y_1!} p_1^{y_1} (1 - p_1)^{(M_{10}-y_1)} \frac{M_{01}!}{(M_{01}-y_0)!y_0!} p_0^{y_0} (1 - p_0)^{(M_{01}-y_0)}$$

5 y_1 is the number of subjects with $y_1 = 1, E = 1$ and $N = 0$; M_{10} is the number with $E = 1$ and $N = 0$; y_0
 6 and M_{01} are the corresponding numbers where $E = 0$ and $N = 1$; the “data” are $y_1, M_{10}, y_0, N_{01}, E$ and N .

7 Prior Distributions for parameters R_{10}, R_{01} , and κ_1

8

9 We use a log-normal prior for κ_1 :

10 3A) $f_{\kappa_1}(\ln(e_1)) = \frac{1}{e_1 \sigma_1 \sqrt{2\pi}} e^{-\frac{1}{2}(\ln(e_1) - \mu_1)^2 / \sigma_1^2}$ for $0 \leq \kappa_1 < \infty$ and 0 elsewhere,

11 where μ_1 and σ_1^2 are the mean and variance of $\ln(e_1)$ The median of e_1 is e^{μ_1} and the variance is
 12 $(e^{\sigma_1^2} - 1)e^{2\mu_1 + \sigma_1^2}$. To double the variance of e_1 (on the original, non-log scale), we solve: $(e^{\sigma_2^2} -$
 13 $1)e^{2\mu_1 + \sigma_2^2} = 2(e^{\sigma_1^2} - 1)e^{2\mu_1 + \sigma_1^2}$ for σ_2 . For $\sigma_1 = 0.54$ (the initial value of in σ_1 used in the Example,
 14 main text), we use $\sigma_2 = 0.67$ to double the variance of e_1 .

15 We use a beta prior for R_{10} :

16 4A) $f_{R_{10}}(p_1) = \frac{\Gamma(\alpha_1 + \beta_1)}{\Gamma(\alpha_1)\Gamma(\beta_1)} p_1^{\alpha_1 - 1} (1 - p_1)^{\beta_1 - 1}$, $0 < p_1 < 1$

17 and a beta prior for R_{01} :

18 5A) $f_{R_{01}}(p_2) = \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0)\Gamma(\beta_0)} p_0^{\alpha_0 - 1} (1 - p_0)^{\beta_0 - 1}$, $0 < p_0 < 1$

19 $f_{R_{10}}(p_1)$ and $f_{R_{01}}(p_0)$ are 0 for p_1 or $p_0 \notin [0, 1]$.

20 The priors $f_{R_{01}}(p_0)$ and $f_{R_{10}}(p_1)$ are non-informative uniform priors if we use $\alpha_j = \beta_j = 1$.

21 Posterior Distribution

22 For analysis of CE_{10} , the posterior distribution is:

23 6A) $f_{R_{10}, R_{01}, \kappa_1 | Y_1, Y_0}(p_1, p_0, e_1 | y_1, y_0)$
 24
$$= \frac{1}{G} p_1^{y_1} (1 - p_1)^{(M_{10}-y_1)} \frac{\Gamma(\alpha_1 + \beta_1)}{\Gamma(\alpha_1)\Gamma(\beta_1)} p_1^{\alpha_1 - 1} (1 - p_1)^{(\beta_1 - 1)} p_0^{y_0} (1 - p_0)^{M_{01} - y_0} \frac{\Gamma(\alpha_0 + \beta_0)}{\Gamma(\alpha_0)\Gamma(\beta_0)}$$

 25
$$\times p_0^{\alpha_0 - 1} (1 - p_0)^{\beta_0 - 1} \frac{1}{e_1} e^{-\frac{1}{2}(\ln(e_1) - \mu_1)^2 / \sigma_1^2}$$

$$\begin{aligned}
 &= \frac{1}{G} p_1^{y_1 + \alpha_1 - 1} (1 - p_1)^{M_{10} - y_1 + \beta_1 - 1} p_0^{y_0 + \alpha_0 - 1} (1 - p_0)^{M_{01} - y_0 + \beta_0 - 1} \frac{1}{e_1} e^{-\frac{1}{2}(\ln(e_1) - \mu_1)^2 / \sigma_1^2} \\
 &\quad \text{for } 0 < p_1 < 1, 0 < p_0 < 1, \text{ and } 0 < e_1 < \infty, \text{ and} \\
 &= 0 \text{ elsewhere,}
 \end{aligned}$$

$$\text{where } \frac{1}{G} = \frac{\Gamma(M_{10} + \alpha_1 + \beta_1)}{\Gamma(y_1 + \alpha_1)\Gamma(M_{10} - y_1 + \beta_1)} \frac{\Gamma(M_{01} + \alpha_0 + \beta_0)}{\Gamma(y_0 + \alpha_0)\Gamma(M_{01} - y_0 + \beta_0)} \frac{1}{\sigma_1 \sqrt{2\pi}}.$$

The posterior distribution for analysis of CE_{10} is the product of two independent *beta* distributions and a *log-normal* distribution. It is therefore straightforward to sample from this joint distribution, by independently sampling: p_1 (for R_{12}) from a $beta(y_1 + \alpha_1, M_{10} - y_1 + \beta_1)$, p_0 (for R_{01}) from $beta(y_0 + \alpha_0, M_{01} - y_0 + \beta_0)$, and e_1 (for κ_1) from a $log-normal(\mu_1, \sigma_1^2)$ distribution.

Evaluation

To evaluate the posterior distribution, the supplemental R code (Web Appendix S3) performs the described sampling in 100,000 independent replications. For each sample, it calculates $\frac{1}{\varepsilon_1} \frac{R_{10}(1 - R_{01} + \varepsilon_1 R_{01})}{R_{01}}$ (or $\frac{R_{10}}{R_{01}(\varepsilon_2 - \varepsilon_2 R_{10} + R_{10})}$), equal to the parameter of interest CE_{10} (or CE_{01}). It then calculates desired statistics from the empiric distribution of sampled values (e.g., 2.5 and 97.5 percentiles for the 95% credibility interval for CE_{10} , and the 50th percentile for the median). The program uses $\alpha_i = \beta_i = 1$ as the default parameters for the *beta* distribution, as this choice yields uninformative, uniform priors. The user must input the parameters of the log-normal prior distribution of κ_i (details and further explication in Web Appendices S1-S2).

1 **Table 1** Summary of Parameter Definitions and relationship to Causal Effects

Parameter Symbol	Definition
M_{en}	Number at risk at baseline in cohort, with $E = e, N = n$ for $e, n = 0,1$
R_{en}	Risk during follow-up among those with $E = e, N = n$ for $e, n = 0,1$
R_{enu}	Risk during follow-up among those with $E = e, N = n, U = u$ for $e, n = 0,1$
R_e	Risk during follow-up among those with $E = e$ for $e = 0,1$
RR_e	Risk ratio comparing risk among those with $E = e, N = 1$ to that among those with $E = e, N = 0$ for $e = 0,1$; e.g., $RR_1 = R_{11}/R_{10}$.
$Y(e)$	Counterfactual value of Y if E were set to e
$R_{en}(e')$	Counterfactual risk among those with $E = e, N = n$, if E were set to e' for $e, n, e' = 0,1$.
$R_{enu}(e')$	Counterfactual risk among those with $E = e, N = n, U = u$, if E were set to e' for $e, n, e' = 0,1$.
CE_{en}	Causal effect of E among those with $E = e, N = n$ for $e, n = 0,1$; $CE_{en} = R_{en}(1)/R_{en}(0)$
$\varepsilon_1 (\varkappa_1)^\ddagger$	Bias Parameter for Probabilistic bias of CE_{10} ; Equation 5, Assumption A5c
$\varepsilon_2 (\varkappa_2)^\ddagger$	Bias Parameter for Probabilistic bias of CE_{01} ; Equation 6, Assumption A5d
$\varepsilon_3, \varepsilon_4^{\ddagger, \dagger}$	Bias Parameter for Probabilistic bias Analyses of CE_{11}, CE_{00} ; Assumption A7

2 [‡]bias parameter for probabilistic bias analysis; corresponding parameter for Bayesian analysis in Parentheses

3 [†]These parameters discussed and used for probabilistic bias analyses for CE_{11} and CE_{00} .

4

5 **Table 2** Distribution of Self-harm Episodes (Y), Hormone Use (E), Prior Vaccination (N)

Variable	- Value of the Variable -							
	1+	0	1+	0	1+	0	1+	0
Number self-harm episodes (Y)	1+	0	1+	0	1+	0	1+	0
Hormone use (E ; yes/no)	yes	yes	no	no	yes	yes	no	no
Vaccinated (N ; yes/no)	yes	yes	yes	yes	no	no	no	no
Number with this combination	10	58	30	152	11	116	42	380

6

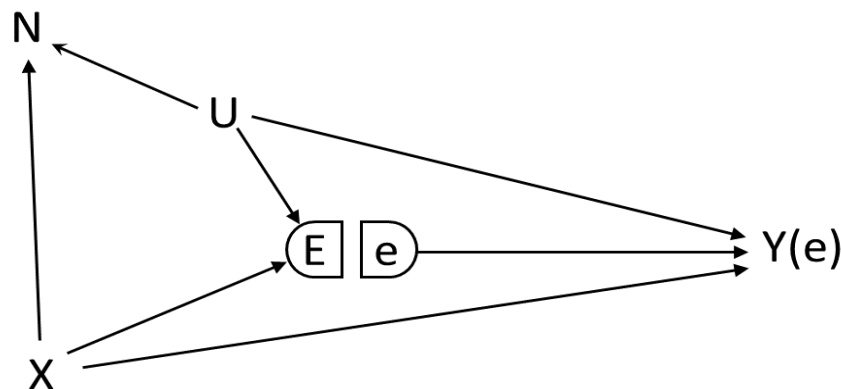
7 **Table 3.** Summary of Estimated Values of Selected, Identifiable Parameters in the Example

Parameter	Estimated Value	Description
cRR	0.88	crude Risk Ratio – association of risk with E
mRR	0.88	Mantel-Haenszel Risk Ratio, adjusted for the negative control
RR_1	1.70	Risk ratio for association of risk with N among those with $E = 1$
RR_0	1.66	Risk ratio for association of risk with N among those with $E = 0$
R_{10}	0.087	Risk among those with $E = 1$ and $N = 0$
R_{01}	0.165	Risk among those with $E = 2$ and $N = 1$

8

1

Figure 1

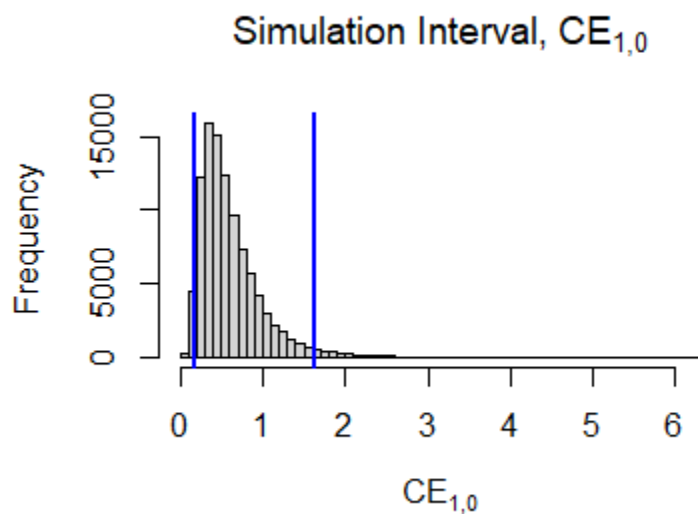


N and Y are expected to be associated, with or without conditioning on E .

2

3 Single World Intervention Template (SWIT)¹⁸ showing causal relationships. E represents the exposure
4 variable, X measured confounder(s), U unmeasured confounder(s), N a negative control and $Y(e)$ the
5 potential outcome of Y if E were set to e . These relationships are assumed correct, and are consistent
6 with assumptions A1–A4.

7



Blue lines indicate 2.5 and 97.5 Percentiles

8

9 Figure 2. Histogram showing the distribution of confounding-adjusted causal effect estimates from
10 probabilistic bias analysis, with the prior distribution of the bias parameter as described for the Example.

11

1 **Online Supplement**

2 **Web Appendix S1**

3 This Web Appendix (S1) includes additional comments about use of the R programs (Web Appendix S2
4 and S3).

5 Probabilistic bias analyses (R program in Appendix S2) require the user to specify a distribution for ε_1 (or
6 ε_2). Specifying a distribution with a median of 1 (e.g., as in the example, main text), would be
7 appropriate if *a priori* information suggested that the counterfactual odds $\frac{R_{10}(0)}{1-R_{10}(0)}$ is expected to exceed
8 the true, identifiable odds $\frac{R_{01}}{1-R_{01}}$ with about 50-50 probability. However, if *a priori* information suggests
9 that the counterfactual odds is more likely than not to be larger than the identifiable odds $\frac{R_{01}}{1-R_{01}}$ or,
10 similarly, that R_{10}/R_{01} is likely to over-estimate the causal effect CE_{10} , then the median should be chosen
11 greater than 1. This might reflect a situation in which exposure was more strongly associated with high-
12 risk values of U than was the negative control. Conversely, if the counterfactual odds is more likely than
13 not to be smaller than the identifiable odds $\frac{R_{01}}{1-R_{01}}$, then the median should be chosen less than 1. The
14 variance reflects uncertainty, with larger variances indicating a wider range of values of ε_1 judged to be
15 compatible with a priori information.

16 Bayesian analyses require the user to specify the prior distribution of κ_1 (corresponding comments also
17 hold for κ_2). To do so using the R program supplied (Appendix S3), the user inputs the median of κ_1 ,
18 which is the mode of the *log-normal* distribution used as the prior for κ_1 . The log-normal parameter μ_1
19 (equation 6A of the Appendix, main text) is the logarithm of this median. The user also inputs the desired
20 value for the ratio of the 10th percentile of κ_1 divided by the median in the distribution of κ_1 . The program
21 uses these inputs to calculate the variance σ_1^2 (equation 6A of the Appendix 2, main text) that yields these
22 user-supplied percentiles. It then generates random samples from the posterior distribution (that depends
23 on the parameters as described. To illustrate, suppose, for example, that the researcher thought that the
24 median was 1 and that the 10th percentile was one-half as large as the median. This relationship would
25 hold if the standard deviation of the log-normal prior were set to 0.54. In this case, the prior would
26 specify that 80% of the values of κ_1 were between 0.5 (the 10 percentile) and 2.0 (the 90 percentile).
27 Since $f_{\kappa_1}(\ln(e_1)|data) = f_{\kappa_1}(\ln(e_1))$, i.e. the posterior and prior are the same for κ_1 , we expect the
28 Bayesian analysis to agree closely with the probabilistic bias analysis.

29 Note S1: The analysis described for CE_{10} (or CE_{01}) is fully Bayesian as all parameters needed to write the
30 (conditional) probability of the data used have a prior, a reflection of the relationships in equations 2A-

1 6A, Appendix 2 of the main text. In particular, a prior distribution for the causal effects CE_{10} (or CE_{01}) is
2 indirectly specified through the priors for R_{10} , R_{01} and κ_1 (or κ_2).

3 Note S2: To measure causal effects using risk odds ratios (rather than risk ratios as above), then we can
4 define $CE_{10}' = \frac{R_{10}(1)/(1-R_{10}(1))}{R_{10}(0)/(1-R_{10}(0))}$. The definitions (of R_{10} , R_{01} , κ_1 and CE_{10}) and distributions above are
5 unchanged (programs to conduct probabilistic bias and Bayesian analyses using risk ratios are given in
6 Web Appendices S2 and S3).

7 Comparison of the prior distribution for κ_1 , illustrated by empirically simulating it (Figure S1) for the
8 situation described in the example (main text), with the corresponding posterior distribution (Figure S2)
9 shows a meaningful shift (e.g. median changes from 0.99 to 0.48), illustrating that the data do provide
10 information about the causal effect under the assumed priors.

11 Note S3: Specification of priors for R_{10} , R_{01} and κ_1 implies a prior for CE_{10} because of the deterministic
12 relationships (equation 5, main text). In view of these inter-relationships of the parameters, specifying a
13 prior for the causal effects themselves, e.g. for CE_{10} , in addition to priors for R_{10} , R_{01} and κ_1 could lead
14 to inconsistencies in view of the deterministic relationships. However, one could specify priors for, say,
15 R_{10} , κ_1 and CE_{10} , but not, say, R_{01} . In this case the posterior distribution would be different from that in
16 Appendix 2 (main text), but could be evaluated using Gibbs sampling or other Markov Chain Monte
17 Carlo approach. Here, the approach implemented in the program specifies a prior for the causal effect
18 (CE_{10}) indirectly through those for R_{10} , R_{01} and κ_1 (Figure S1; the prior for example of main text), and
19 does not separately specify a prior for CE_{10} .

20 Note S4: If direct evidence, *a priori* considerations or otherwise expected differences between the E and
21 N conditional distributions (departures from assumption A5) suggest that $R_{10}(0)$ will tend to exceed R_{01}
22 so that κ_1 is less than 1 (in Assumption A5a), that would suggest that E tends to more strongly associate
23 with high-risk values of U than does N . This pattern suggests that $R_{01}(1)$ should tend to be less than R_{10}
24 and that κ_2 may be greater than 1. These considerations suggest, that if the prior distribution of κ_1 has
25 substantial mass below 1, then that for κ_2 might reasonably be specified to have substantial mass above 1
26 (and conversely).

27 **Web Appendix S2**

28 Appendix S2 is an R program for Probabilistic Bias Analysis; effects are measured using risk ratios (See
29 also Web Appendix S1).

```
30 # Probabilistic Bias Analysis for Causal Effect CE21 (or CE12), using Risk Ratios  
31 library(tidyverse)  
32 library(pscl)  
33 library(dplyr)
```

```
1 library(data.table)
2 library(readxl)
3 # ----- user inputs ----- #
4 med = 1 #user input, desired median (log(med)=mean yields mu, of log-normal dist)
5 q10 = 0.5 #user input, desired ratio: 10-percentile/med: used to calculate sd
6 n=100000 #user input, number of random draws
7 N = 1
8
9 ## ----- input data -----
10 dat <- array( rep(0,8), dim=c(2,2,2))
11 dat[1,1,1]= 10 #P(Y=1, E=1, NC=1)
12 dat[2,1,1]= 58
13 dat[1,2,1]= 30
14 dat[2,2,1]= 152
15
16 dat[1,1,2]= 11 #P(Y=1,E=1,NC=2)
17 dat[2,1,2]= 116
18 dat[1,2,2]= 42
19 dat[2,2,2]= 380
20
21 n1 = dat[1,1,2]+dat[2,1,2]
22 n2 = dat[1,2,1]+dat[2,2,1]
23
24 ## ----- END: input data -----
25
26 save = NULL
27
28 # ----- calculate conditional Probs from input dat, so they to sum to 1
29 for (j in 1:2){
30   for (k in 1:2){
31     dat[1,j,k]=dat[1,j,k]/sum(dat[,j,k])
32     dat[2,j,k]=1-dat[1,j,k]
33   }
34 }
35 # ----- END: calculate cond probs.
36
37 AVERAG<-0
38 use=NULL
39 # =====
40 # ===== Crude Risk Ratio =====
41 RR1= dat[1,1,1]/dat[1,1,2]
42 RR2= dat[1,2,1]/dat[1,2,2]
43 cRR = (dat[1,1,1]+dat[1,1,2])/sum(dat[,1,])
44 cRR = cRR/((dat[1,2,1]+dat[1,2,2])/sum(dat[,2,]) )
45 R11= dat[1,1,1]
46 R12= dat[1,1,2]
47 R21= dat[1,2,1]
48 R22= dat[1,2,2]
49
50 mRR= dat[1,1,1]*(dat[1,2,1]+dat[2,2,1])/sum(dat[,1])
51 mRR=mRR+ dat[1,1,2]*(dat[1,2,2]+dat[2,2,2])/sum(dat[,2])
```

```
1 den= dat[1,2,1]*(dat[1,1,1]+dat[2,1,1])/sum(dat[,1])
2 den=den+ dat[1,2,2]*(dat[1,1,2]+dat[2,1,2])/sum(dat[,2])
3 mRR= mRR/den
4 mRR
5 #RD_E.NC #association of E, NC
6
7 crude=data.frame(RR1=RR1, RR2=RR2, R11=R11, R12=R12, R21=R21, R22=R22, cRR) #RD for
8 Exposure, then for NC
9 crude
10 c(R11/R12, R21/R22) #Y-assoc with N|E=1 or 2
11 #
12 # ----- start log-normal, prior for e1 -----;
13 z10=qnorm(0.10, mean=0, sd = 1) #get z-score @ 10-percentile
14 sd = (log(q10)-log(med))/z10 #set sd on log scale so (log(q10)-log(med))/sd)=z10
15 z90=qnorm(0.90, mean=0, sd = 1)
16 c(z10, z90)
17 c(exp( log(med)+z10*sd),exp( log(med)-z10*sd) )
18 e1= rnorm(n, log(med), sd)
19 #e2= rnorm(n, -log(med), sd) #use symmetric value for e2, per Note S4
20 e1=exp(e1)
21 c(quantile(e1, probs=c(0.025, 0.5, 0.975)), sd) #display 95% interval for e1
22 #e2=exp(e2)
23 #c(exp(qnorm(0.10, mean=0, sd = sd) ),
24 # exp(qnorm(0.90, mean=0, sd = sd) ) ) #display 20%, 80% of epsilon
25 #reverse the order
26 R12_0 = (e1*R21/(1-R21))/(1+(e1*R21/(1-R21))) # calc R12(0) from R21 1st (use obs'd R21)
27 #r21 = rbinom(n, n2, R21)/n2 # sample R21, with random error
28 r12_0 = rbinom(n, n2, R12_0)/n2 # sample R12(0), with random error
29 r12 = rbinom(n, n1, R12 )/n1 # sample R12, with random error
30 # in Example, v. similar result if change, sample R21 & R12 1st, correct w/ e1 last
31 # CE12= r12*(1-r21+e1*R21)/(e1*r21) # use this if change order, RR for CE10; modify for ROR
32 CE12 = r12/r12_0 # RR for CE10
33 #CE21= r12/(r21*(e2-e2*r12+r12)) # use this if desire RR; modify for ROR
34 # ----- Figure 2 -----
35 hist(CE12, breaks=50, main=expression("Simulation Interval, CE"[paste(1,",",0)]),
36 xlab = expression("CE"[paste(1,",",0)]),
37 sub="Blue lines indicate 2,5 and 97.5 Percentiles", cex.sub = 0.8)
38 abline(v=quantile(CE12, probs= 0.975),col="blue",lwd=2)
39 abline(v=quantile(CE12, probs= 0.025),col="blue",lwd=2)
40 quantile(CE12, probs=c(0.025, 0.5, 0.975)) # display empiric 95% credibility interval
41
42 Web Appendix S3
43 R program for Bayesian analyses, described in Appendix 2, main text (see also Web Appendix S1);
44 effects measured using risk ratios.
45 # Bayesian Analysis for Causal Effect CE12 (or CE21), using Risk Ratios
46 library(tidyverse)
47 library(pscl)
48 library(dplyr)
49 library(data.table)
50 library(readxl)
```

```
1 # ----- user inputs ----- #
2 med = 1 #user input, desired median (log(med)=mean yields mu, of log-normal dist)
3 ratio_10 = 0.5 #user input, desired ratio: 10-percentile/med: used to calculate sd
4 n=100000 #user input, number of random draws
5
6 ## ----- input data -----
7 dat <- array( rep(0,8), dim=c(2,2,2))
8 dat[1,1,1]= 10 #P(Y=1, E=1, NC=1)
9 dat[2,1,1]= 58
10 dat[1,2,1]= 30
11 dat[2,2,1]= 152
12
13 dat[1,1,2]= 11 #P(Y=1,E=1,NC=2)
14 dat[2,1,2]= 116
15 dat[1,2,2]= 42
16 dat[2,2,2]= 380
17 # ===== ----- Some Risks, ORs, & RRs ----- =====
18 OR1= (dat[1,1,1]/dat[2,1,1])/(dat[1,1,2]/dat[2,1,2])
19 OR2= (dat[1,2,1]/dat[2,2,1])/(dat[1,2,2]/dat[2,2,2])
20
21 RR1= (dat[1,1,1]/dat[1,1,2])/((dat[1,1,1]+dat[2,1,1])/(dat[1,1,2]+dat[2,1,2]))
22 RR2= (dat[1,2,1]/dat[1,2,2])/((dat[1,2,1]+dat[2,2,1])/(dat[1,2,2]+dat[2,2,2]))
23 R11= dat[1,1,1]
24 R12= dat[1,1,2]/(dat[1,1,2]+dat[2,1,2])
25 R21= dat[1,2,1]/(dat[1,2,1]+dat[2,2,1])
26 R22= dat[1,2,2]
27 cRR = (dat[1,1,1]+dat[1,1,2])/((dat[1,1,1]+dat[2,1,1]+dat[1,1,2]+dat[2,1,2]))
28 cRR = cRR /((dat[1,2,1]+dat[1,2,2])/((dat[1,2,1]+dat[2,2,1]+dat[1,2,2]+dat[2,2,2])))
29 aRR = dat[1,1,1]*(dat[1,2,1]+dat[2,2,1]) + dat[1,1,2]*(dat[1,2,2]+dat[2,2,2])
30 aRR=aRR/( dat[1,2,1]*(dat[1,1,1]+dat[2,1,1]) + dat[1,2,2]*(dat[1,1,2]+dat[2,1,2]))
31
32 crude=data.frame(OR1=OR1, OR2=OR2, RR1=RR1, RR2=RR2, R11=R11, R12=R12, R21=R21,
33 R22=R22) #RD for Exposure, then for NC
34 (R12/(1-R12))/(R21/(1-R21))
35 crude
36 y1= dat[1,1,2]
37 n1= dat[2,1,2]+y1
38 y2= dat[1,2,1]
39 n2= dat[2,2,1]+y2
40 # -- first, sample r21 & r12
41 r21= rbeta(n, y2+1, n2-y2+1)
42 r12= rbeta(n, y1+1, n1-y1+1)
43 # -- second, sample e1 from log-normal prior
44 # use input median and ratio of median to 10th percentile: get sd
45 z10=qnorm(0.10, mean=0, sd = 1) #get z-score @ 10-percentile
46 sd = (log(ratio_10)-log(med))/z10 #set sd on log scale so (log(ratio_10)-log(med))/sd=z10
47 z90=qnorm(0.90, mean=0, ratio_10, sd = 1)
48 c(z10, z90, med, sd)
49 c(exp( log(med)+z10*sd),(med),exp( log(med)-z10*sd) ) #display e1 at 10%, 50%, 90%
50 #sd= sd3 #<----- set sd=sd3 for sensitivity analyses to double Var of prior
51 #if sd= sd3 NOT commented out, doubles Var (value from end of pgm)
```

```
1 #for Example, sd=0.674 <-- use to double variance from initial value of 0.54
2 e1= rnorm(n, log(med), sd) #sd, sd3 in a sens Anal w/ double Var
3 c(exp(qnorm(0.10, mean=log(med), sd = sd) ),
4   exp(qnorm(0.90, mean=log(med), sd = sd) ) ) #display empiric 20%, 80% of e1
5 e1=exp(e1)
6 # ----- END: sample e1 from log-normal prior -----;
7
8 #quantile( e1, probs=c(0.025, 0.5, 0.975)) #display 2.5th 97.5 percentiles
9 CE12= r12*(1-r21+e1*R21)/(e1*r21) # use this for RR, could do ROR
10
11 #CE21= r12/(r21*(e2-e2*r12+r12)), for RR
12 # ----- Figure S2 -----
13 hist(CE12, breaks=50, main=expression("Posterior Distribution CE"[paste(1,",",0)]),
14     xlab = expression("CE"[paste(1,",",0)]),
15     sub="Blue lines indicate 95% Credibility Interval", cex.sub = 0.8)
16 abline(v=quantile(CE12, probs= 0.975),col="blue",lwd=2)
17 abline(v=quantile(CE12, probs= 0.025),col="blue",lwd=2)
18 quantile(CE12, probs=c(0.025, 0.5, 0.975)) # display empiric 95% credibility interval
19
20 c(R12/(R12+dat[2,1,2]), R21/(R21+dat[2,2,1]), RR1, RR2, cRR, aRR)
21 # ----- empirically evaluate prior for CE -----
22 prior=FALSE
23 if (prior==TRUE){
24   r21= rbeta(n, 1, 1)
25   r12= rbeta(n, 1, 1)
26   e1= rnorm(n, log(med), sd)
27   e1=exp(e1)
28   s2=sd**2
29   c(var(e1), (exp(s2) -1)*exp(s2) )
30   CE12= r12*(1-r21+e1*R21)/(e1*r21) # use this for RR, could do ROR
31   quantile(CE12, probs=c(0.025, 0.5, 0.975))
32   parm2 = as_tibble(CE12, CE=CE12)
33   names(parm2)="CE"
34   parm2
35   parm2 = subset(parm2, parm2$CE< 30)
36   # ----- Figure S1 -----
37   hist(parm2$CE, breaks=50, main=c(expression("Prior Distribution CE"[paste("1",",",",0")])*~"(truncated
38   to display)")),
39     sub="Blue lines indicate 2.5th and 97.5 percentile", cex.sub = 0.8,
40     xlab = expression("CE"[paste(1,",",0)]))
41   abline(v=quantile(parm2$CE, probs= 0.975),col="blue",lwd=2)
42   abline(v=quantile(parm2$CE, probs= 0.025),col="blue",lwd=2)
43   }
44 # for use in Sensitivity (Bayesian analyses) find sd* that double var on e-scale
45 #solver - arbitrary function called fx(x), one variable
46 fx<- function(x){ #just an example
47   s1=0.54**2 #var used (log-scale), solve for x for Var doubles (real scale)
48   2*(exp(2*s1)-exp(s1) ) - (exp(2*x)-exp(x) ) }
49 sd3=sqrt(uniroot(fx,interval=c(-10,3))$root )
50 c(sd, sd3 ) #set sd=sd3 for sensitivity analyses to double Var of prior
```


1 **Web Appendix S4**

2 Lipsitch et al¹ note that an ideal negative control exposure would have the same causes in common with
 3 the outcome as did the actual exposure. Figure S3 is a single world intervention graph in which E and N
 4 have different causes in common with Y , but wherein N can still serve as a negative control exposure.
 5 Assumption (A5), while still possible, may be now less plausible.

6 **Web Appendix S5**

7 Claim 1: Under assumption (A1-A4) and (A5a) or (A5b): (i) $R_{10}(0) = R_{01}$; and (ii) CE_{10} is identified by
 8 the ratio of observable risks R_{10}/R_{01} .

9 Proof: Proofs of (i) under assumption A5a and (ii) are in the main text. Here we show that (A1-A4) and
 10 (A5b) implies (i). By conditional independence of E and N (Assumption A1) and rules of conditional
 11 probabilities:

12 1S) $p(U = u|E = 1, N = 0) = p(E = 1|U = u)p(N = 0|U = u)p(U = u)$, for all u , and:

13 2S) $p(U = u|E = 0, N = 1) = p(E = 0|U = u)p(N = 1|U = u)p(U = u)$, for all u .

14 From the last line of Equation 4) of the main text:

$$\begin{aligned}
 15 \quad R_{10}(0) &= \frac{\sum_u R_{01u} P(E=1|U=u)P(N=0|U=u)p(U=u)}{\sum_u p(E=1|U=u)p(N=0|U=u)p(U=u)} && \text{(Equation 4, main text)} \\
 16 \quad &= \sum_u R_{01u} p(U = u|E = 1, N = 0) && \text{(by Equation 1S)} \\
 17 \quad &= \sum_u R_{01u} p(U = u|E = 0, N = 1) && \text{(by Assumption A5b)} \\
 18 \quad &= \frac{\sum_u R_{01u} P(E=0|U=u)P(N=1|U=u)p(U=u)}{\sum_u p(E=0|U=u)p(N=1|U=u)p(U=u)} && \text{(by Equation 2S)} \\
 19 \quad &= R_{01} && \text{(equals } R_{01}, \text{ a weighted average).}
 \end{aligned}$$

20 Note: Equations 1S and 2S show that Assumption (A5b) is implied by (A5a) and that (A5b) is weaker.

21
 22 Claim 2: Under assumptions (A1-A4, A5b), $CE_{10} = CE_{01}$.

23 Proof: The proof parallels that of Claim 1. $R_{en}(e')$ is the weighted average of the counterfactual
 24 outcomes $R_{enu}(e')$, weighted by $P(U = u|E = e, N = n) = P(E = e|U = u)P(N = n|U = u)p(U =$
 25 $u)/\sum_u P(E = e|U = u)P(N = n|U = u)p(U = u)$, so:

$$\begin{aligned}
 26 \quad CE_{01} &= \frac{R_{01}(1)}{R_{01}(0)} \\
 27 \quad &= \frac{\sum_u R_{01u}(1)P(E=0|U=u)P(N=1|U=u)p(U=u)}{\sum_u p(E=0|U=u)p(N=1|U=u)p(U=u)} / \frac{\sum_u R_{01u}(0)P(E=0|U=u)P(N=1|U=u)p(U=u)}{\sum_u p(E=0|U=u)p(N=1|U=u)p(U=u)} \\
 28 \quad &= \sum_u R_{01u}(1)P(U = u|E = 0, N = 1) / \sum_u R_{01u}(0)P(U = u|E = 0, N = 1) \\
 29 \quad &= \sum_u R_{10u}(1)P(U = u|E = 0, N = 1) / \sum_u R_{10u}(0)P(U = u|E = 0, N = 1) \\
 30 \quad &= \sum_u R_{10u}(1)P(U = u|E = 1, N = 0) / \sum_u R_{10u}(0)P(U = u|E = 1, N = 0)
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{\sum_u R_{10u}(1)P(E=1|U=u)P(N=0|U=u)p(U=u)}{\sum_u p(E=1|U=u)p(N=0|U=u)p(U=u)} / \frac{\sum_u R_{10u}(0)P(E=1|U=u)P(N=0|U=u)p(U=u)}{\sum_u p(E=1|U=u)p(N=0|U=u)p(U=u)} \\
 &= \frac{R_{10}(1)}{R_{10}(0)} = CE_{10}
 \end{aligned}$$

The second equality follows by properties of the counterfactuals $R_{01}(1)$ and $R_{01}(0)$, the third equality by equations 1S and 2S, the fourth by assumptions (A1) and (A2), the fifth by assumption (A5b), the sixth by equations 1S and 2S, and the last two by definitions.

Web Appendix S6

This Web Appendix contains additional, basic descriptive information about the cohort in the Example of the main text. First, we note that, during the follow-up period, the age of cohort members with gender-affirming hormone or puberty suppression therapy, was very similar to the age of those without that therapy by definition of the cohort. In particular, the therapy had to have been received by age 20, and the follow-up for the outcome of interest was from age 20 to age 21. Thus, age during the time at risk should not have led to meaningful confounding. Table S1 includes additional descriptive information for cohort of the example in the main text and shows reasonable degree of balance between treatment groups – so these factors should also not have led to meaningful confounding. Consistent with this, we did not adjust for additional covariates and did no modeling or smoothing; causal effect estimates for the example can be calculated from just the data in Table 2 (main text).

Table S1. Descriptive Information for Cohort of Example in the Main Text ^a

Characteristic	Gender-Affirming Therapy ^b N (%)	No Gender-Affirming Therapy N (%)
Ethnicity		
Hispanic	171 (28.3)	44 (22.6)
Non-Hispanic Black	58 (9.6)	13 (6.7)
Other, including Non-Hispanic White	375 (62.1)	138 (70.8)
Gender Identity		
Trans-feminine	261 (43.2)	89 (45.6)
Trans-masculine	322 (53.3)	106 (54.4)
Transgender with unclear sex assigned at birth	21 (3.5)	0 (0)
Age at Index Date ^c (Years)		
Mean (SD)	16.6 (1.9)	17.2 (1.8)
All Subjects	604 (75.6)	195 (24.4)

^a Percentages may not sum to 100 due to rounding

^b Receiving gender-affirming hormones at or before age 20

^c Index date refers to the day of the first evidence of transgender status in Kaiser Permanente EMR

18

19

20

21

22

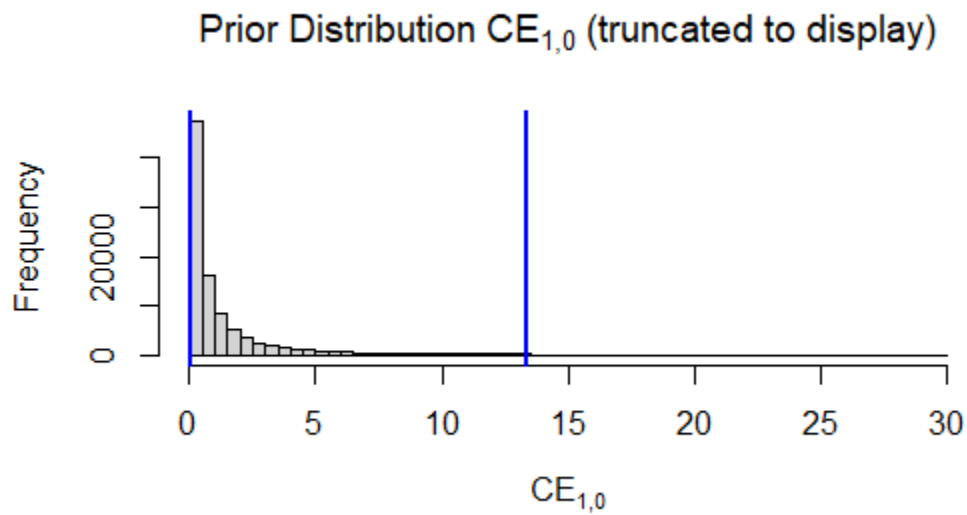
1

2 **References used in Web material**

- 3 1. Lipsitch M, Tchetgen E, Cohen T. Negative Controls: A Tool for Detecting Confounding and Bias
4 in Observational Studies. *Epidemiology*. 2010;21(3):383-388.

5

1 Web Appendix, Figure S1

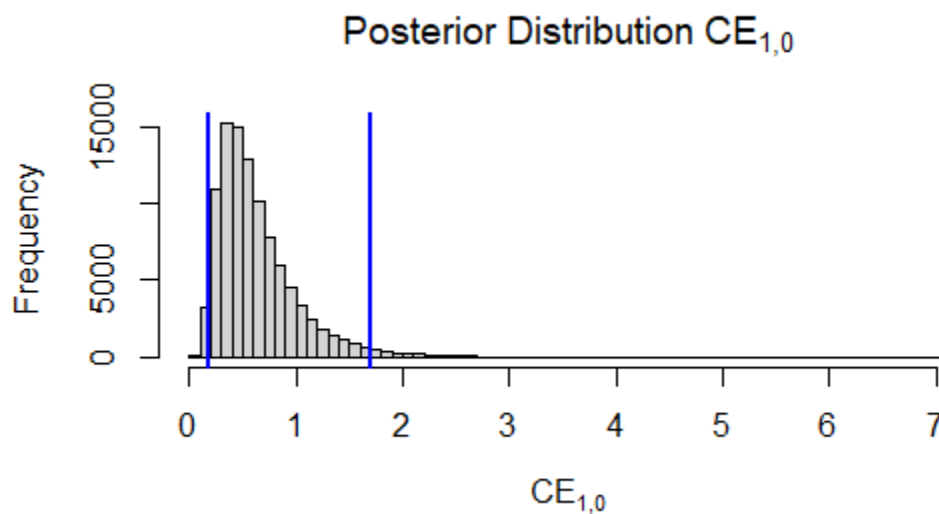


2 Blue lines indicate 2.5th and 97.5 percentile

3 Figure S1 describes the prior distribution of $CE_{2,1}$ for the example in the main text; see also Web
4 Appendix S1. Values of $CE_{2,1}$ were truncated at 30 to display.

5

1 Web Appendix, Figure S2



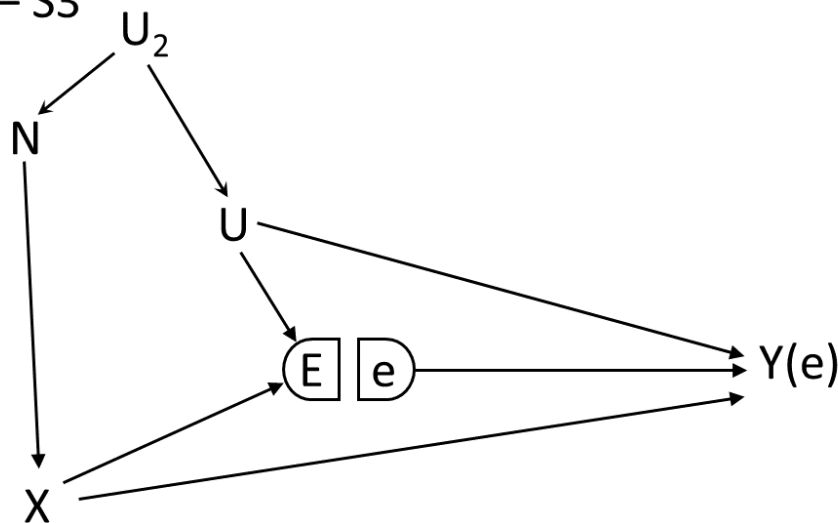
2 Blue lines indicate 95% Credibility Interval

3 Figure S2 describes the posterior distribution of $CE_{2,1}$ for the example in the main text; see also Web
4 Appendix S1. (Values not truncated.)

5

6 Web Appendix, Figure S3

Figure – S3



E and N have different causes in common with Y, but N can still serve as a negative control exposure; the equi-distributional assumption (A5a) or (A5b) may be less plausible.

7