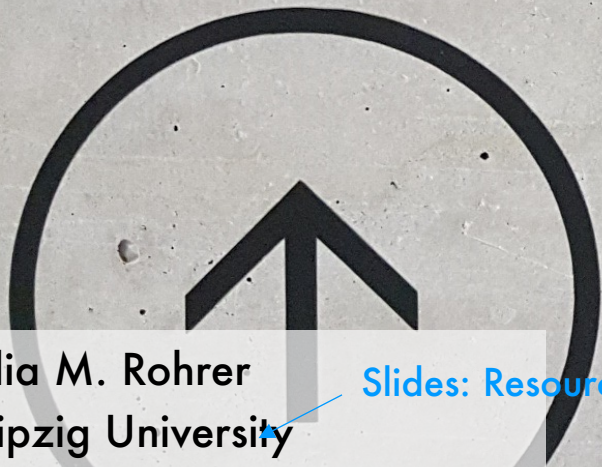



# An introduction to causal graphs & Causal graphs for missing data



Julia M. Rohrer

Leipzig University

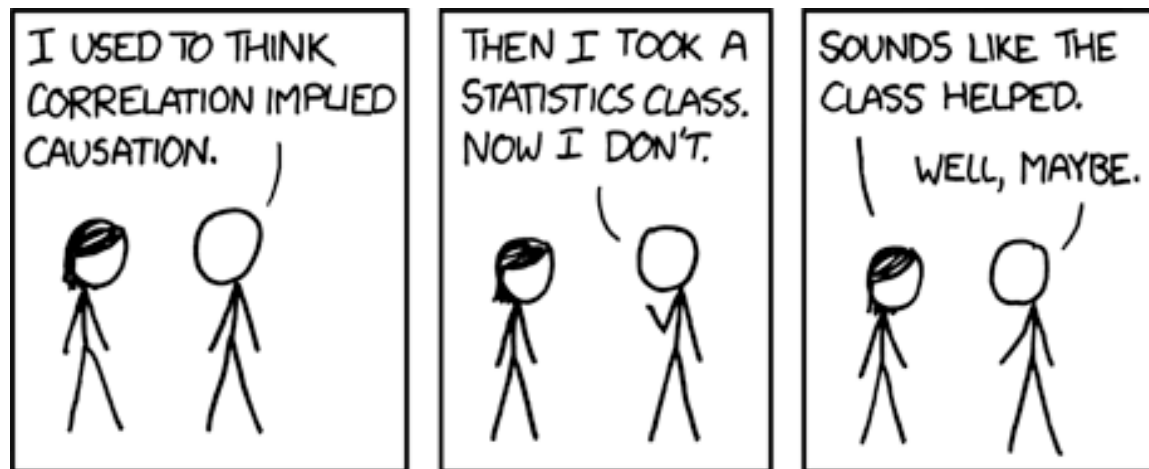
[www.juliarohrer.com](http://www.juliarohrer.com)

 [@dingdingpeng.the100.ci](https://twitter.com/dingdingpeng.the100.ci)

[www.the100.ci](http://www.the100.ci)

[Slides: Resources](#)

# Causal Inference 101



» Option 1: Run an experiment

» Option 2: Give up

» include covariates, maybe?

» or maybe use longitudinal data???

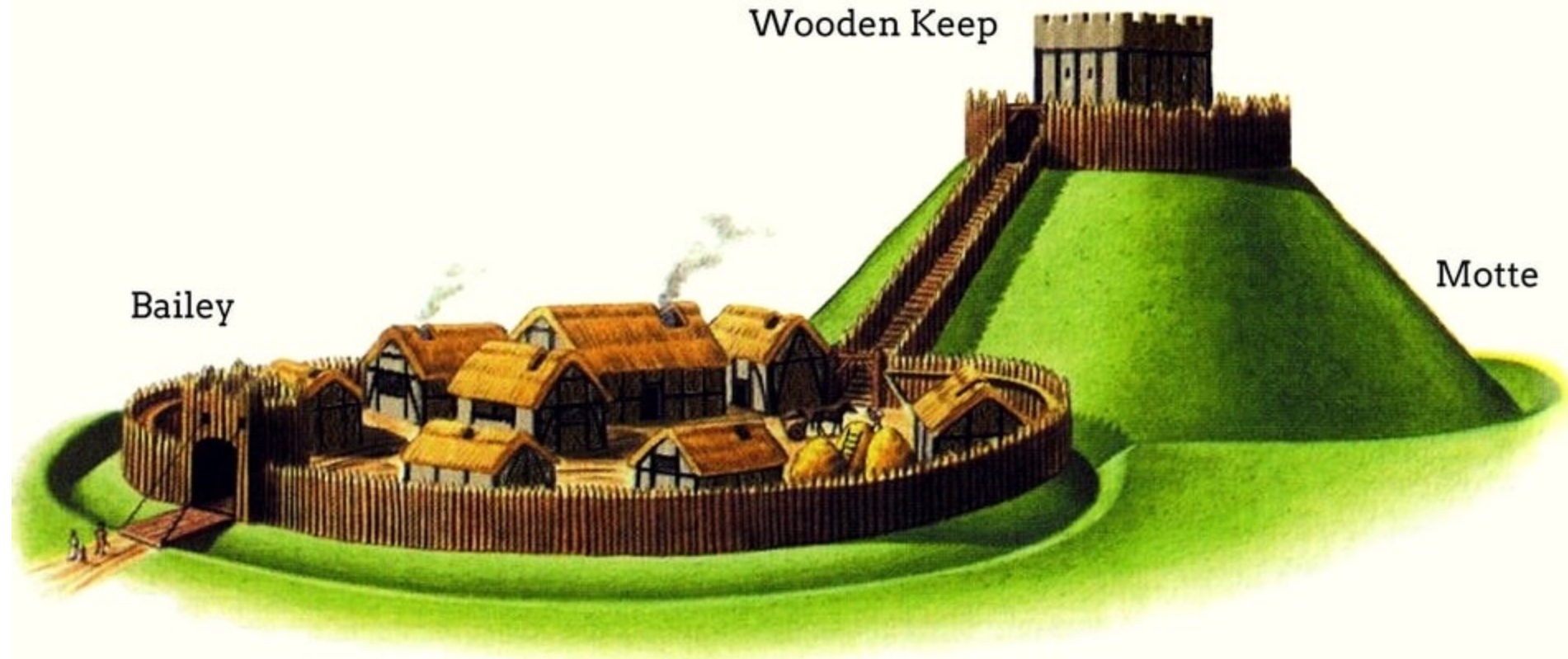
# It's bad for non-experimental research

1. virtually all interesting research questions concern causality
  2. observational data is not admissible for causal inference
1. + 2. = ???

# The non-experimental psych workaround

- » introduction: relies on a causal reading of the literature
- » methods & results: implicitly (but never explicitly) causal inference-y
  - » X predicts Y...even after accounting for...
  - » X is a risk factor for Y
  - » longitudinal associations
- » discussion: only makes sense in terms of causality
- » **IMPORTANT:** add a paragraph that your study was only observational and no causal conclusions are warranted
  - » future longitudinal or experimental studies will surely fix this problem

# Causal inference Motte-and-Bailey



X has a causal effect on Y.

X „predicts“ Y.

# It's bad for experimental researchers



» Causal inference issues that remain despite randomization

» **Missing data**

» Mediation analysis

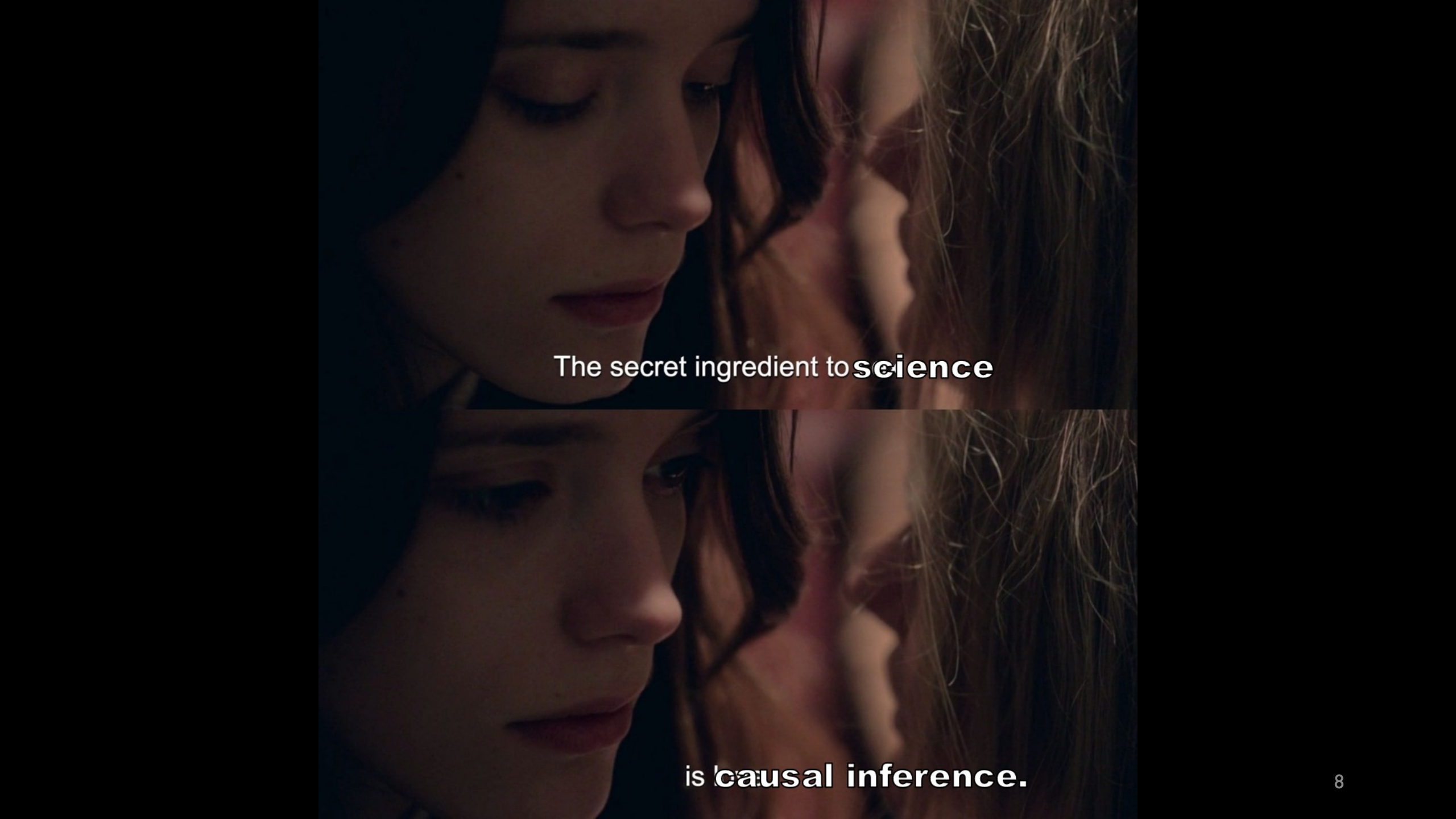
» Moderation analysis

» Generalization across populations, settings, independent and dependent variables

» Measurement

» ...





The secret ingredient to **science**

is **causal inference.**

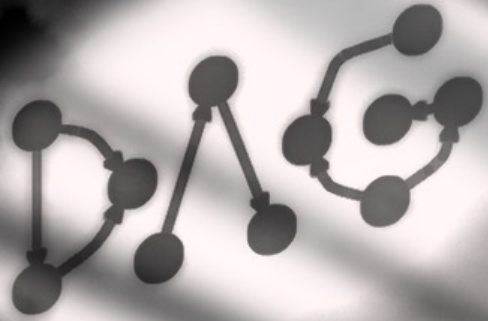
# Causal inference issues

## Identification

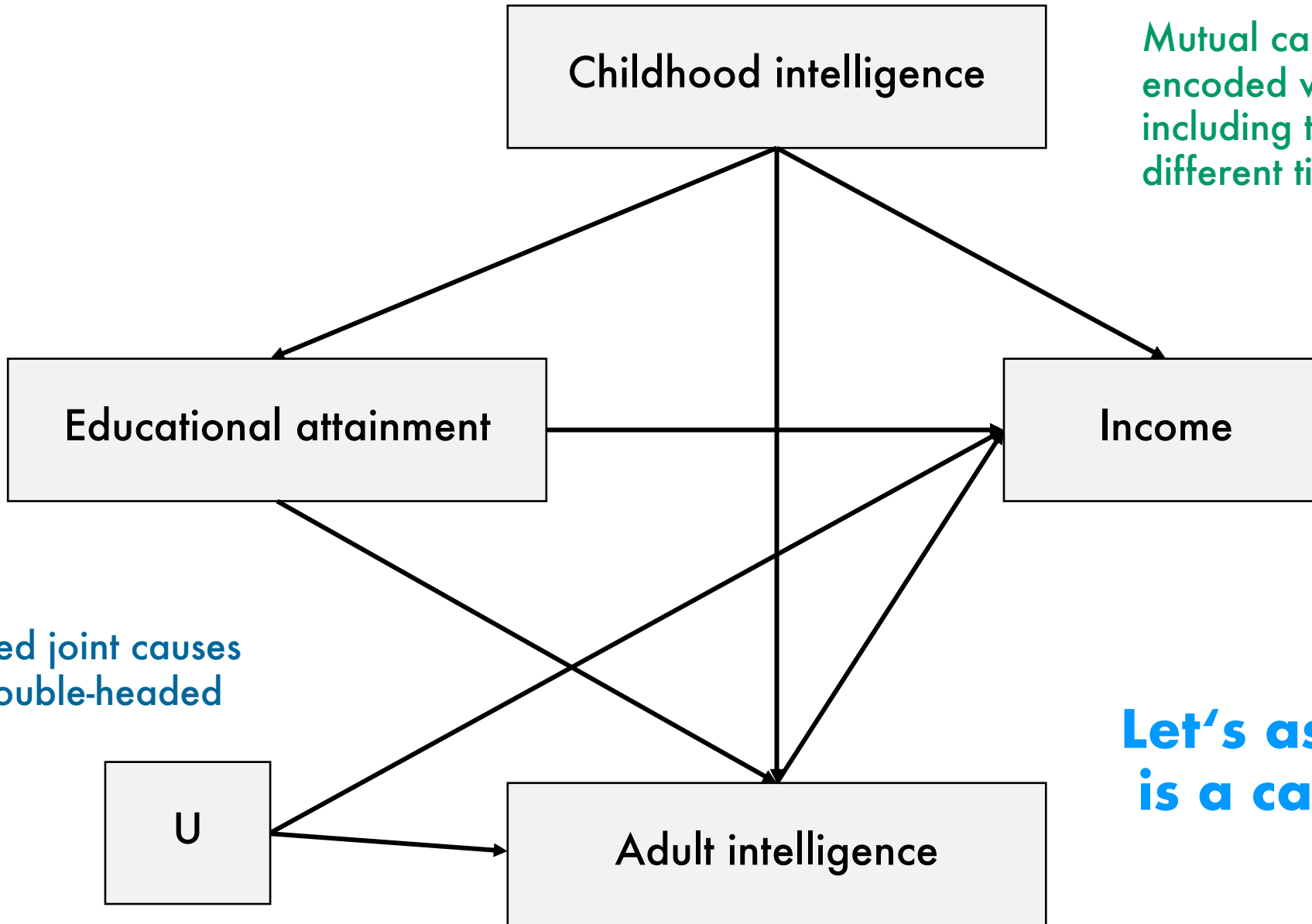
Given your assumptions, if your sample was infinitely large—would you be able to estimate the causal effect of interest?

## Estimation

Actually estimating the effect with the data you got + your assumptions



# Directed Acyclic Graphs

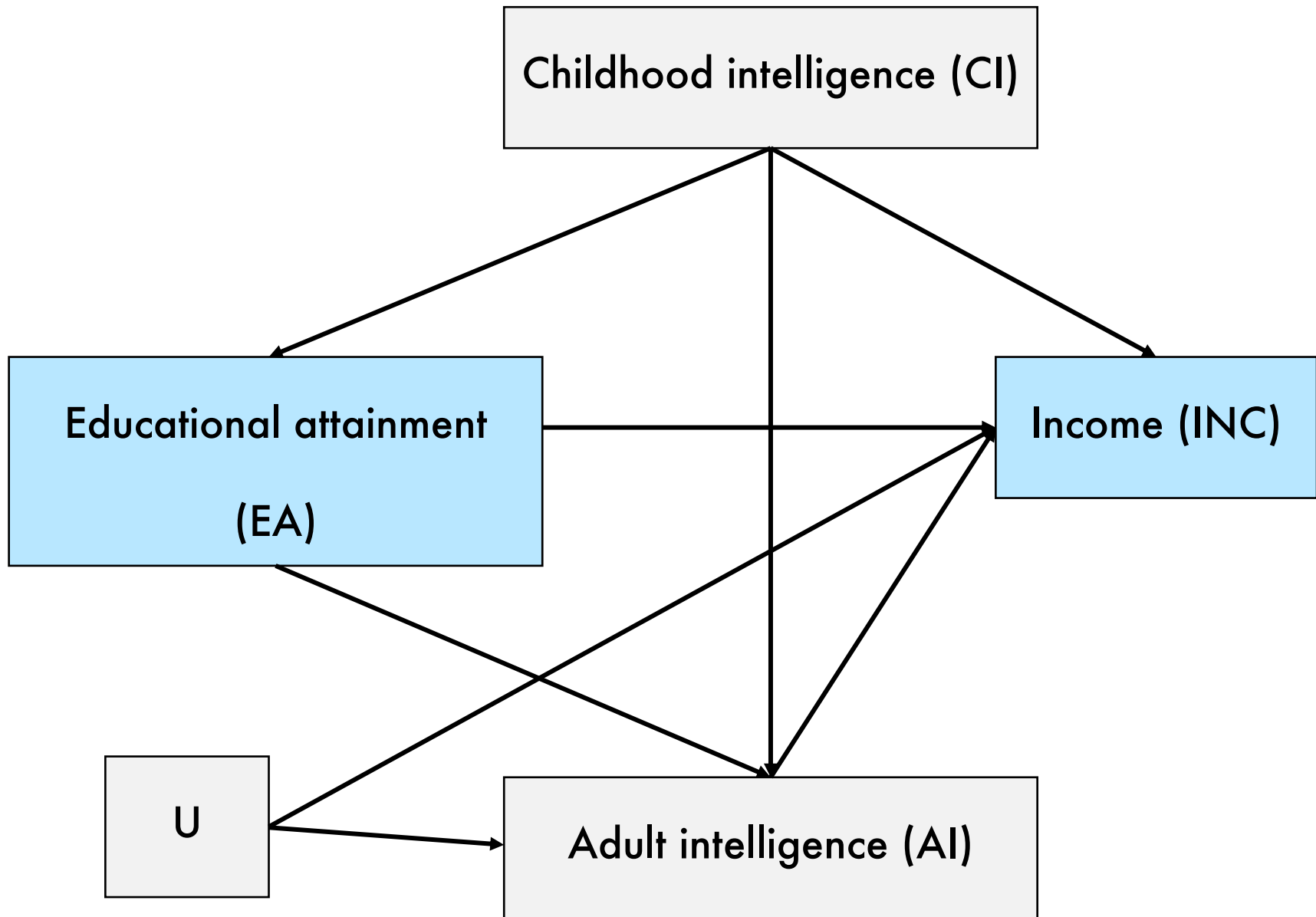


Mutual causation can be encoded without cycles by including the same variable at different time points

Unobserved joint causes replace double-headed arrows

**Let's assume this is a causal DAG**

# Backdoor Criterion



EA → INC  
 EA ← CI → INC

EA ← CI → AI → INC  
 EA ← CI → AI ← U → INC

EA → AI → CI → INC  
 EA → AI → INC  
 EA → AI ← U → INC

Which of these paths lead to associations between EA and INC?

Are these associations causal or non-causal?

→ keep the causal ones, block the non-causal ones

# 3 Fundamental structures

» Chains:  $X \rightarrow M \rightarrow Y$

» transmits a causal association from  $X$  to  $Y$

» conditioning on (control for)  $M$  blocks transmission

# What does it mean to condition on a variable?

- » Any means of introducing information about the variable into the analysis (using it as a control variable, a covariate...)
- » For example
  - » Statistical adjustment in an ANCOVA, regression, SEM...
  - » Using the variable for matching, weighting
  - » Stratification by the variable
  - » Only including participants with a certain level of a variable

# 3 Fundamental structures

» Chains:  $X \rightarrow M \rightarrow Y$

» transmits a causal association from  $X$  to  $Y$

» conditioning on (control for)  $M$  blocks transmission

» if you're interested in the (total) effect of  $X$  on  $Y$ , you usually don't want to do that, because the chain is part of the causal effect of interest  $\rightarrow$  control for mediator leads to overcontrol bias

» aka mediation

# 3 Fundamental structures

» Forks:  $X \leftarrow C \rightarrow Y$

» induces a non-causal association between X and Y

» conditioning on C removes the non-causal association

» you usually want to get rid of these, so conditioning on confounders is usually the way to go

» aka confounding

# Inverted forks $X \rightarrow Z \leftarrow Y$

» does not transmit any association

» e.g., sickness  $\rightarrow$  attendance of this talk  $\leftarrow$  interest in causal inference

» conditioning on  $Z$  opens transmission of a non-causal association

» e.g., conditioning on attendance

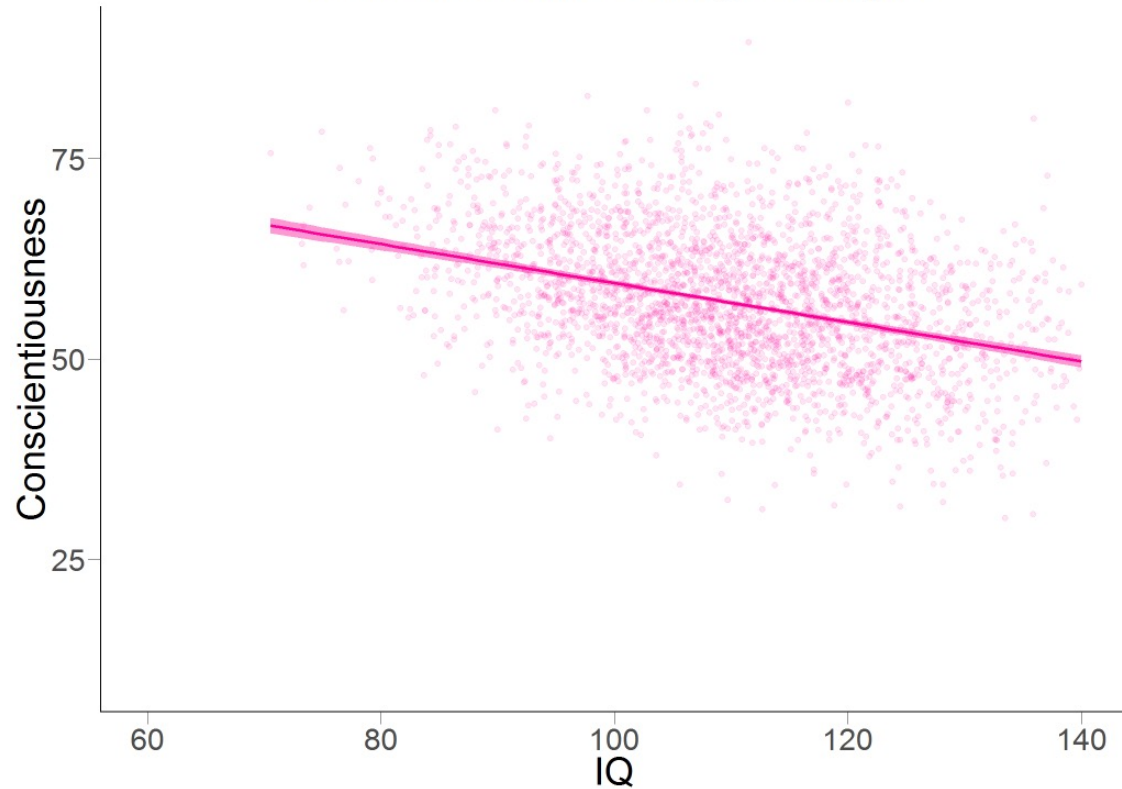
» among the people who attend this talk, those who are sick probably are particularly interested in causal inference

» in other words, third-variable control *introduces* a new spurious association

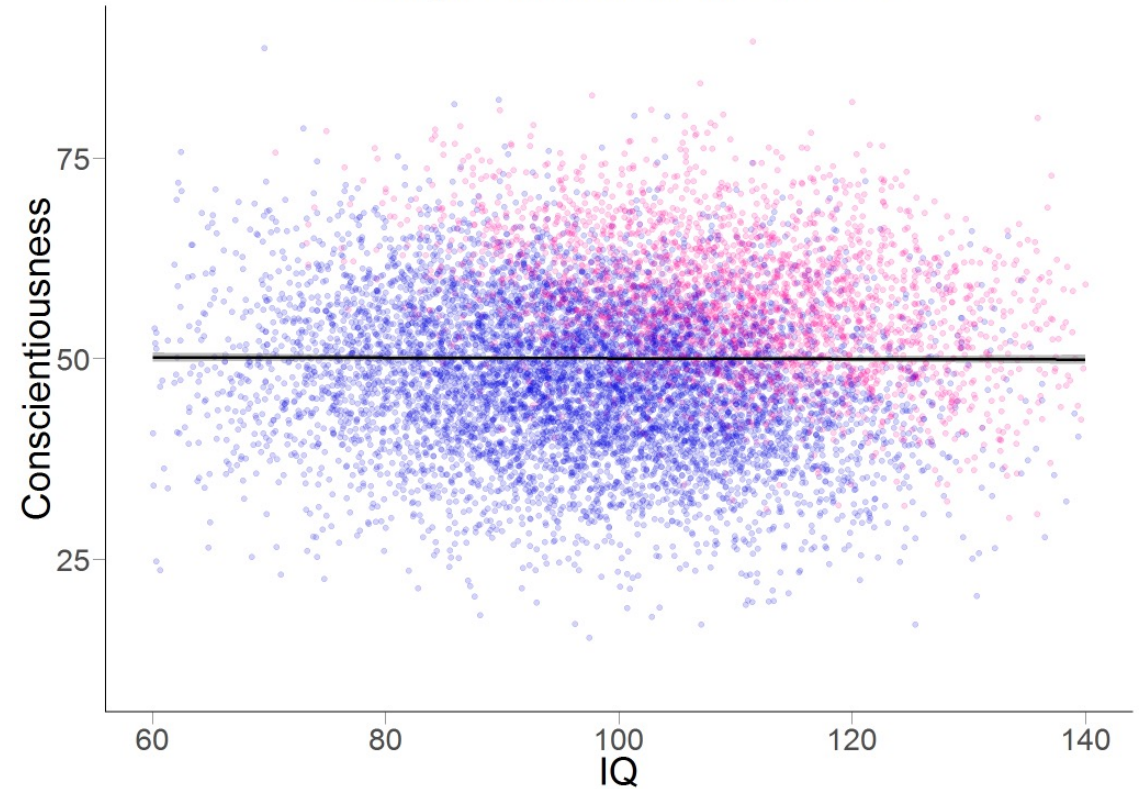
» aka collider bias

# Conscientiousness $\rightarrow$ College $\leftarrow$ IQ

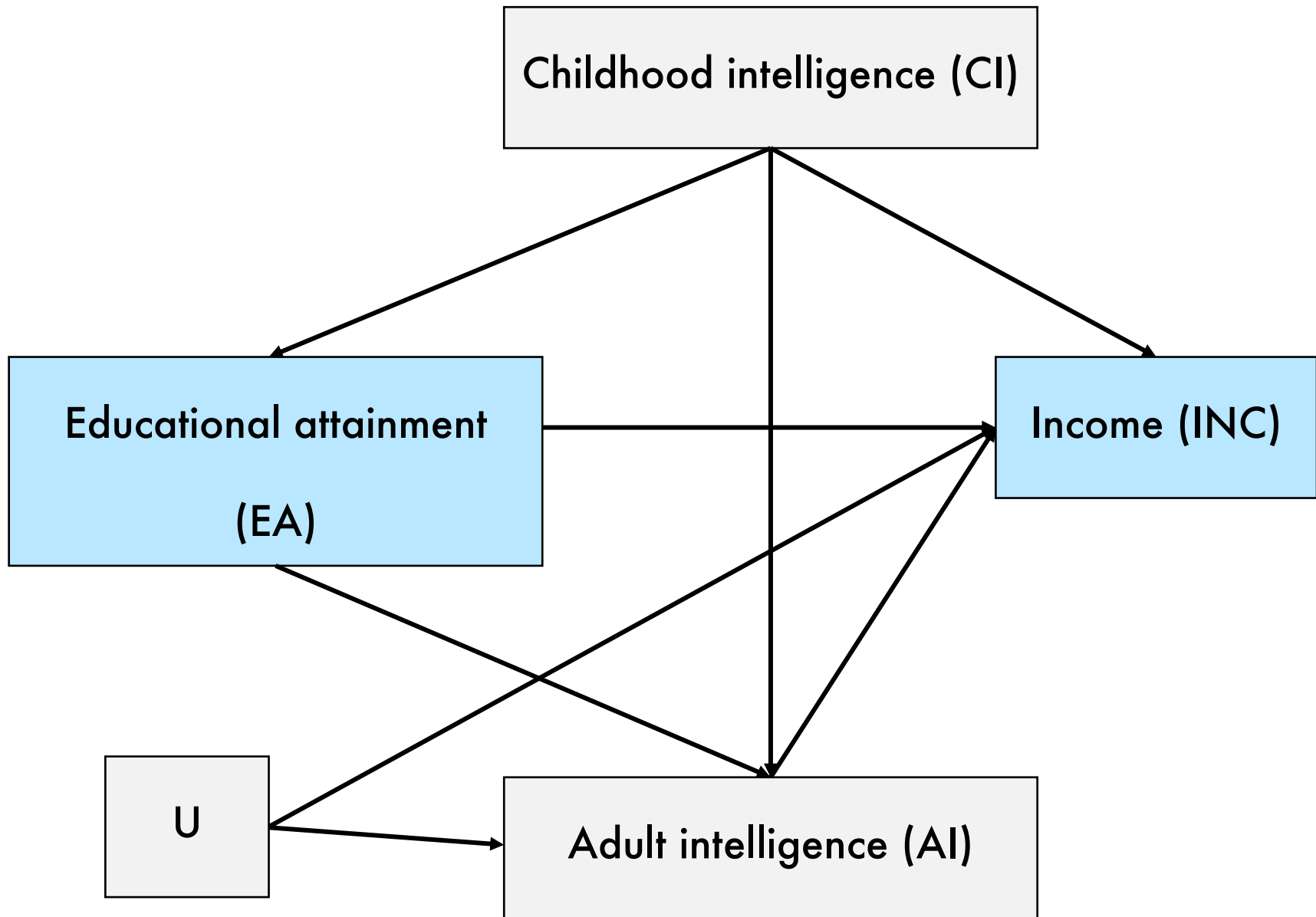
Your awesome college sample



Sample with non-WEIRDos



	Association between X and Y when not conditioning on anything	Association between X and Y when conditioning on the variable in the middle
Chain: $X \rightarrow \text{Mediator} \rightarrow Y$	<b>Causal association</b>	No association (overcontrol bias)
Fork: $X \leftarrow \text{Confounder} \rightarrow Y$	Non-causal association	<b>No association</b>
Inverted fork: $X \rightarrow \text{Collider} \leftarrow Y$	<b>No association</b>	Non-causal association (collider bias)



EA -> INC  
 EA <- CI -> INC

EA <- CI -> AI -> INC  
 EA <- CI -> AI <- U -> INC

EA -> AI -> INC  
 EA -> AI <- U -> INC

# Paths

» Paths that transmit non-causal associations and need to be blocked (backdoor paths)

» EA  $\leftarrow$  CI  $\rightarrow$  INC

» EA  $\leftarrow$  CI  $\rightarrow$  AI  $\rightarrow$  INC

Mediator passes on the non-causal association

Conditioning on the mediator AI would close this one backdoor path (but there are other reasons why we wouldn't want to do that)

Confounder that introduces the non-causal association

Conditioning on the confounder closes both backdoor paths

# Paths

» Paths that transmit causal associations and should not be blocked

» EA -> INC

» EA -> AI -> INC



Arrows all flow from cause to outcome – all fine, don't control away the good stuff

# Paths

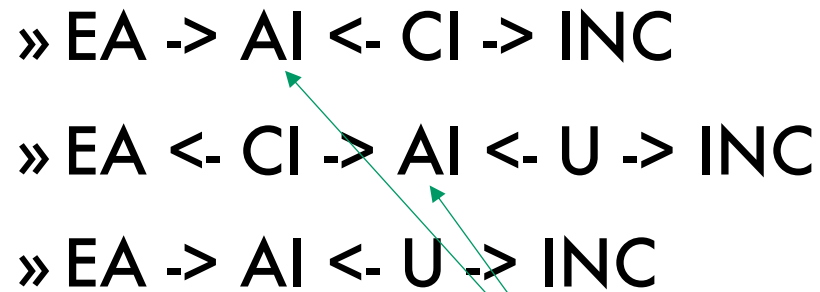
» Paths that are blocked thanks to a collider...

» ...but would lead to non-causal associations *if* the collider was conditioned on

» EA -> AI <- CI -> INC

» EA <- CI -> AI <- U -> INC

» EA -> AI <- U -> INC



Two arrows pointing into node → collider → this path doesn't do anything, you're all good  
UNLESS you condition on the collider

# Paths

» What happens if you condition on the collider?

EA -> AI <- CI -> INC



EA <- CI -> AI <- U -> INC

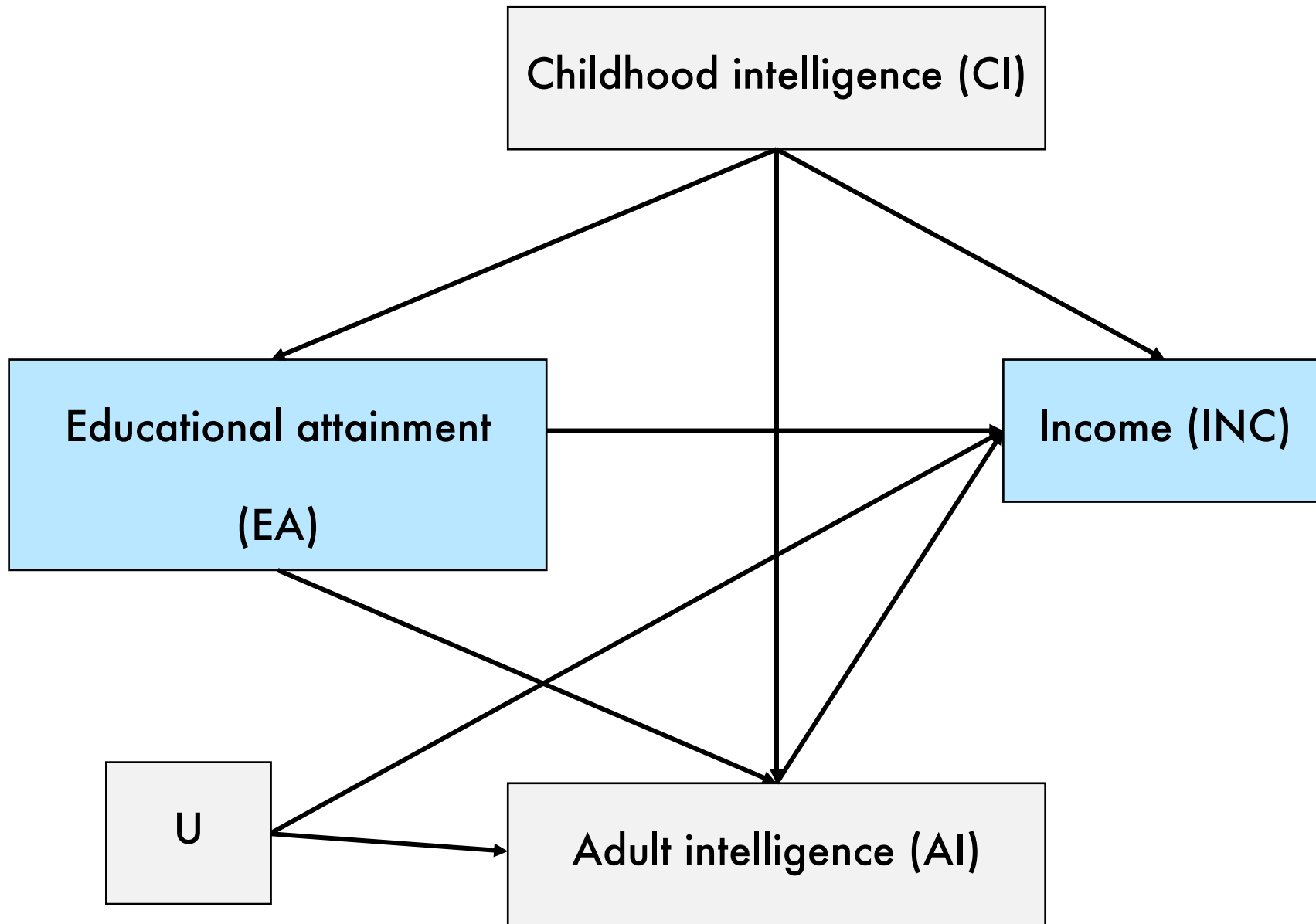


EA -> AI <- U -> INC



New backdoor paths that  
could be fixed by  
conditioning on CI

New backdoor paths that can only be fixed if U  
is measured and conditioned on



*Assuming that this is the correct causal structure:*

*Conditioning on childhood intelligence is both sufficient and necessary to identify the causal effect of interest*

# Causal estimation

- » different ways to condition on a variable, for example
  - » stratification, sub-group analysis
  - » regression adjustment
  - » weighting & matching approaches
- » if your causal identification strategy was wrong, no amount of estimation can rescue you
- » if your causal identification strategy was right, things can still go wrong
  - » Insufficient control due to misspecified models, measurement error in covariates...

*But Wait...*  
**There's**  
**MORE!**

**Causal graphs for missing data**

# Missing data

- » Not all participants always provide all necessary data
- » Often presented as a purely statistical/predictive problem
  - » how can we predict the missing values as well as possible?
- » Somewhat confusing terminology
  - » MCAR, MAR, NMAR (or MNAR)



*Heeeey, MCARena!*

# Missing data

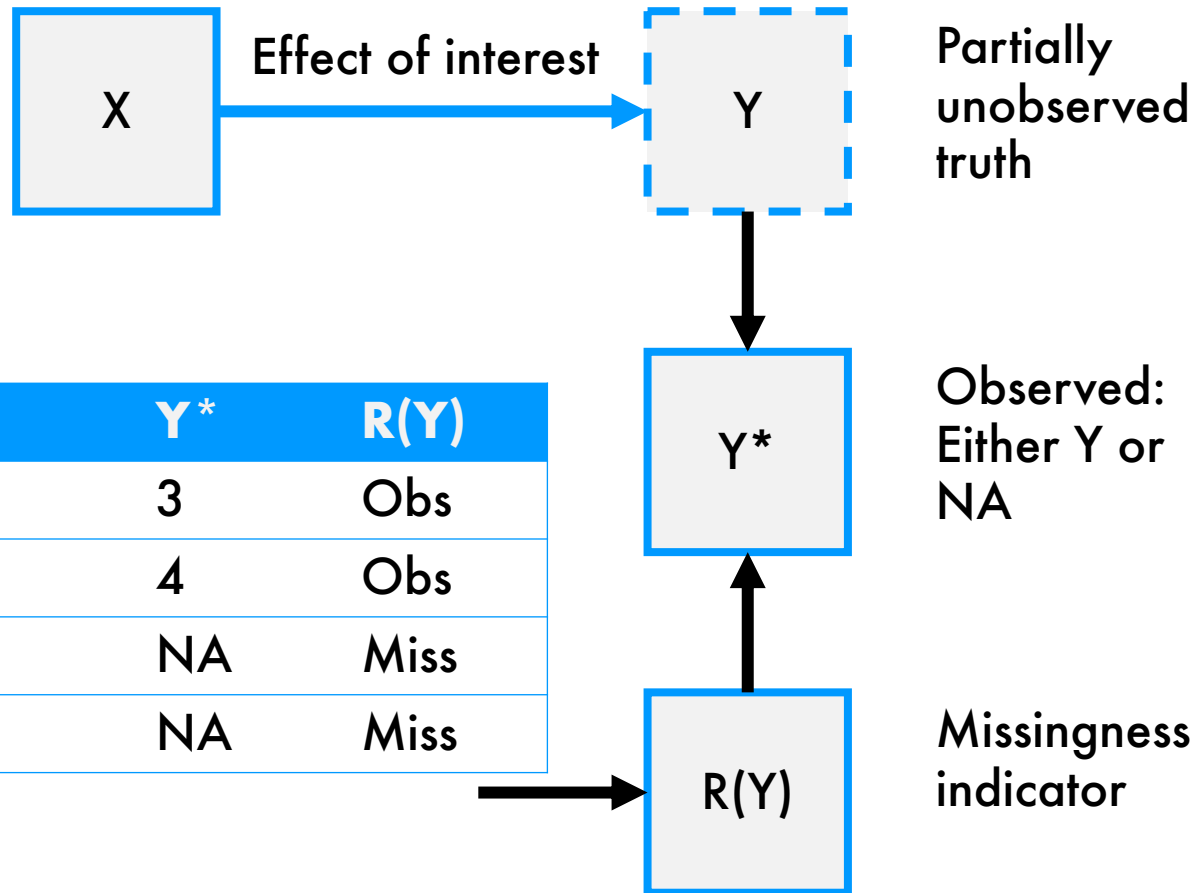
» Causal angle ([Thoemmes & Mohan, 2015](#))

» Focus on *why* the data are missing

» Which variables need to be accounted for to recover an unbiased estimate of some effect of interest

» Sometimes, using a variable that predicts the outcome for imputation purposes can make things worse

# Missing Completely at Random



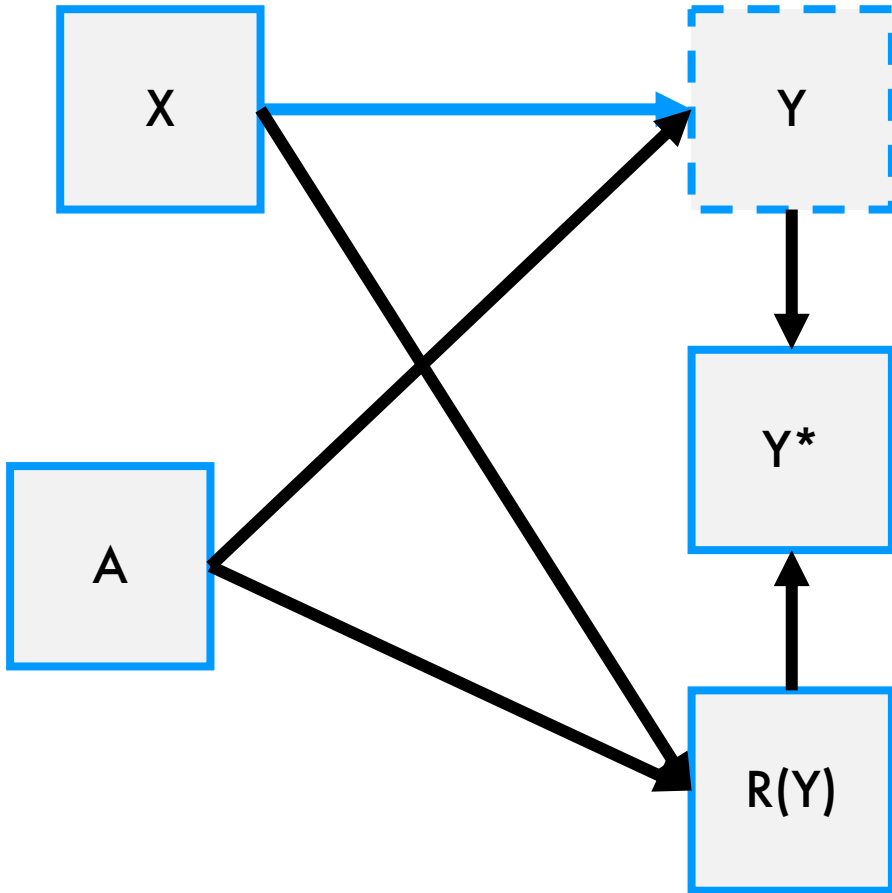
Thoemmes & Mohan (2015)

Crucial question to recover the effect of interest:  
Can we separate Y and R(Y)?

Here:  
 $Y \rightarrow Y^* \leftarrow R(Y)$  blocked because of the collider

Nothing *needs* to be done for an unbiased estimate (complete cases analysis would work)

# Missing at Random



Two open paths between Y and R(Y):

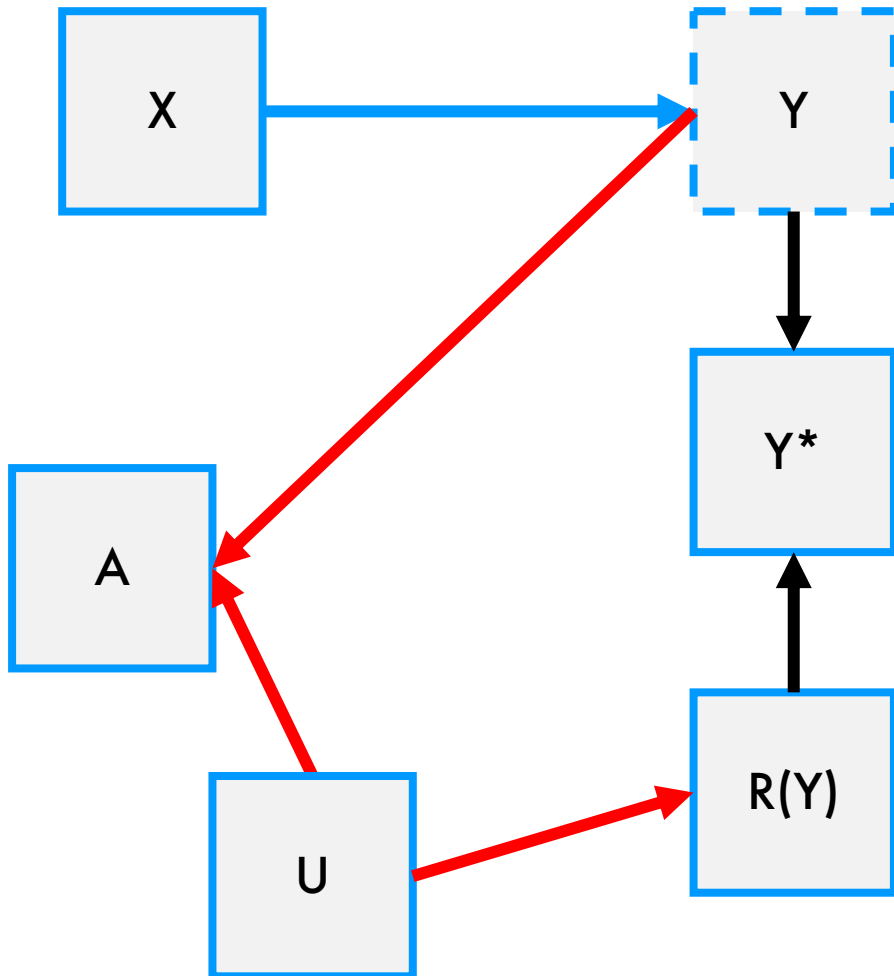
$Y \leftarrow X \rightarrow R(Y)$

$Y \leftarrow A \rightarrow R(Y)$

Block these! e.g.

- multiple imputation
  - FIML estimation
  - inverse probability weighting
- using both X and A

# Missingness is a Causal Inference Problem

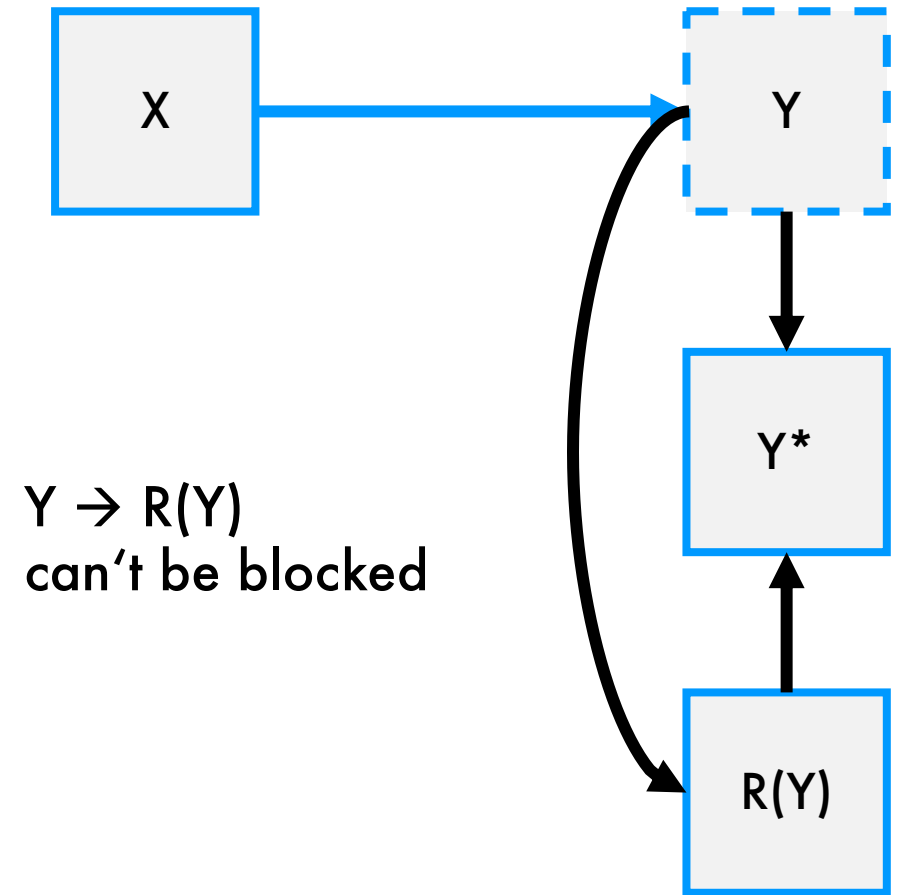
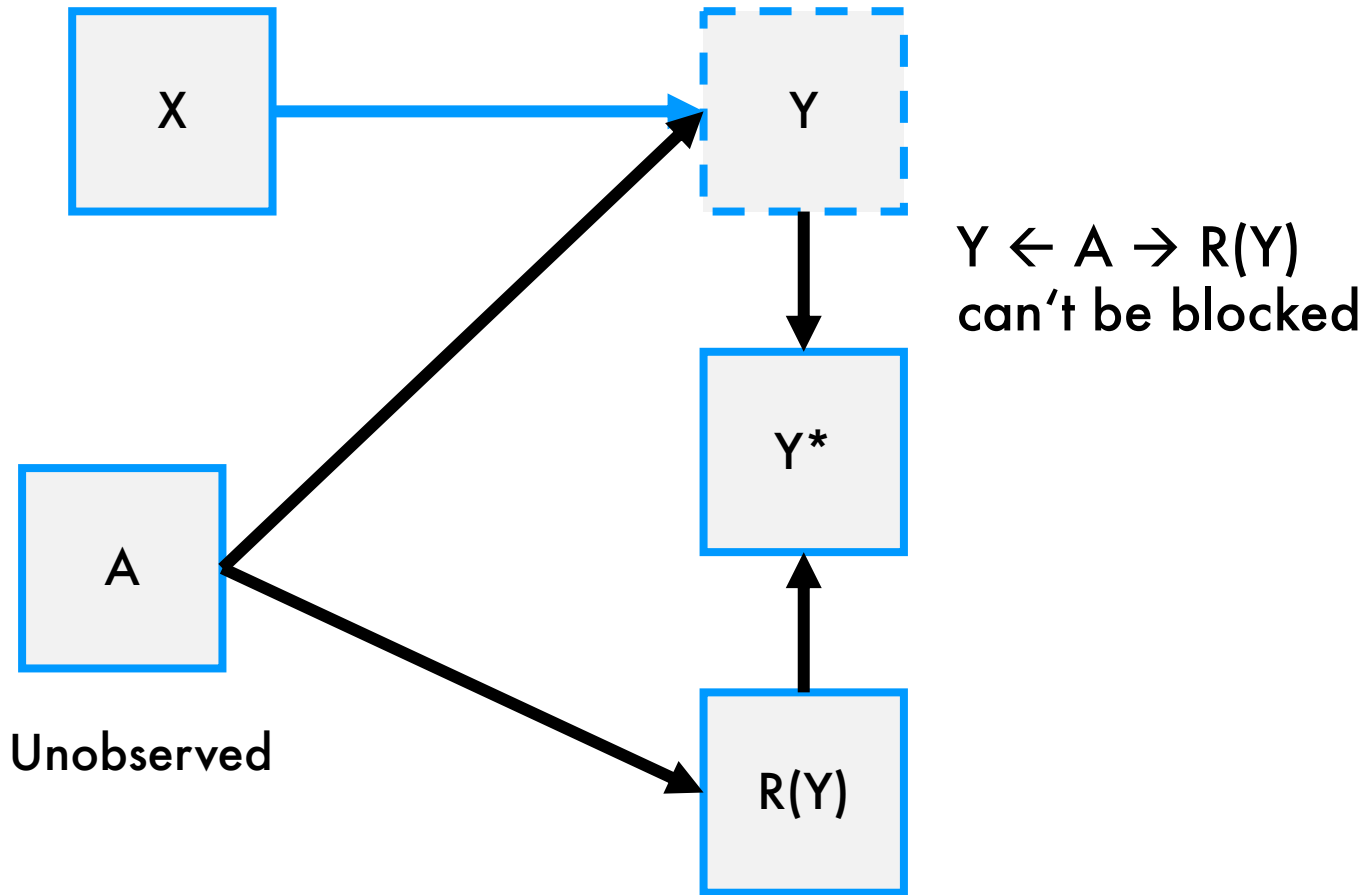


$$Y \rightarrow A \leftarrow U \leftarrow R(Y)$$

is blocked because  $A$  is a collider



Introducing information about  $A$  (e.g., by using it as an auxiliary variable) will *bias* estimates (unless  $U$  is also included)

# Not Missing at Random



# DAGitty — draw and analyze causal diagrams

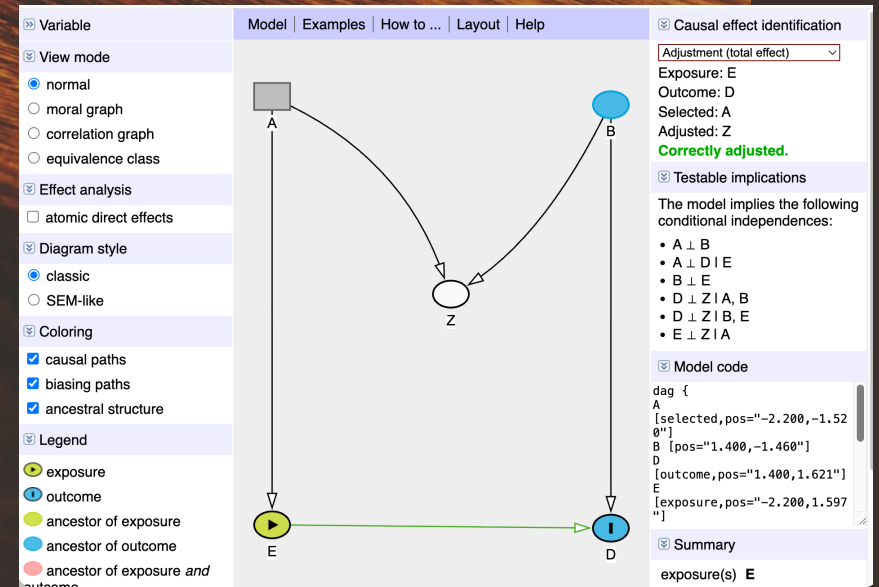
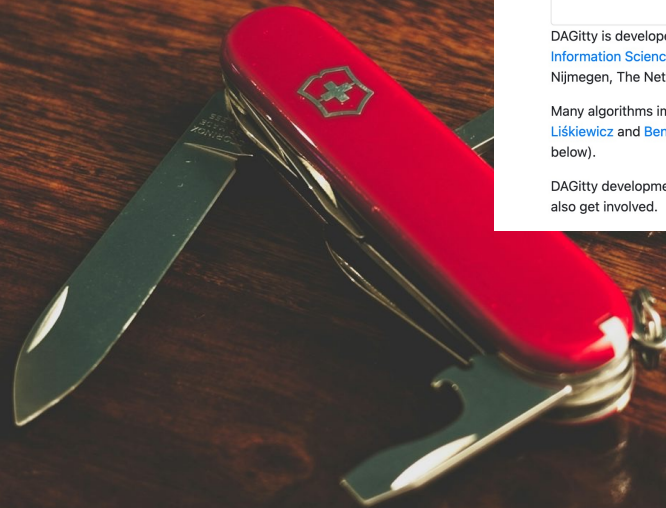
DAGitty is a browser-based environment for creating, editing, and analyzing causal diagrams (also known as directed acyclic graphs or causal Bayesian networks). The focus is on the use of causal diagrams for minimizing bias in empirical studies in epidemiology and other disciplines. For background information, see the ["learn"](#) page.

Launch	Download	Learn	Code
 Launch DAGitty online in your browser.	 Download DAGitty's source for offline use.	 Learn more about DAGs and DAGitty.	 The R package "dagitty" is available on CRAN or github.

DAGitty is developed and maintained by [Johannes Textor](#) (Institute for Computing and Information Sciences, Radboud University, and Medical BioSciences department, Radboudumc, Nijmegen, The Netherlands).

Many algorithms implemented in DAGitty were developed in close collaboration with [Maciej Liśkiewicz](#) and [Benito van der Zander](#), University of Lübeck, Germany (see literature references below).

DAGitty development happens on [github](#). You can download all source code from there and also get involved.



The screenshot displays the DAGitty interface with a causal diagram and analysis results. The diagram shows variables A, B, E, and D, with arrows indicating causal relationships: A → Z, B → Z, E → D, and A → D. The interface includes a sidebar with settings for variable, view mode, effect analysis, diagram style, coloring, and legend. The main panel shows the diagram and analysis results, including causal effect identification and testable implications.

**Variable**

- View mode
  - normal
  - moral graph
  - correlation graph
  - equivalence class
- Effect analysis
  - atomic direct effects
- Diagram style
  - classic
  - SEM-like
- Coloring
  - causal paths
  - biasing paths
  - ancestral structure
- Legend
  - exposure
  - outcome
  - ancestor of exposure
  - ancestor of outcome
  - ancestor of exposure and outcome

**Model** | Examples | How to ... | Layout | Help

**Causal effect identification**

Adjustment (total effect) [v]  
Exposure: E  
Outcome: D  
Selected: A  
Adjusted: Z  
**Correctly adjusted.**

**Testable implications**

The model implies the following conditional independences:

- $A \perp B$
- $A \perp D | E$
- $B \perp E$
- $D \perp Z | A, B$
- $D \perp Z | B, E$
- $E \perp Z | A$

**Model code**

```
dag {
  A
  [selected, pos="-2.200, -1.52
  0"]
  B [pos="1.400, -1.460"]
  D
  [outcome, pos="1.400, 1.621"]
  E
  [exposure, pos="-2.200, 1.597
  "]
}
```

**Summary**

exposure(s) E

# Thank you for your attention!

Julia M. Rohrer

Leipzig University

[www.juliarohrer.com](http://www.juliarohrer.com)

 [@dingdingpeng.the100.ci](https://dingdingpeng.the100.ci)

[www.the100.ci](http://www.the100.ci)

[Slides: Resources](#)