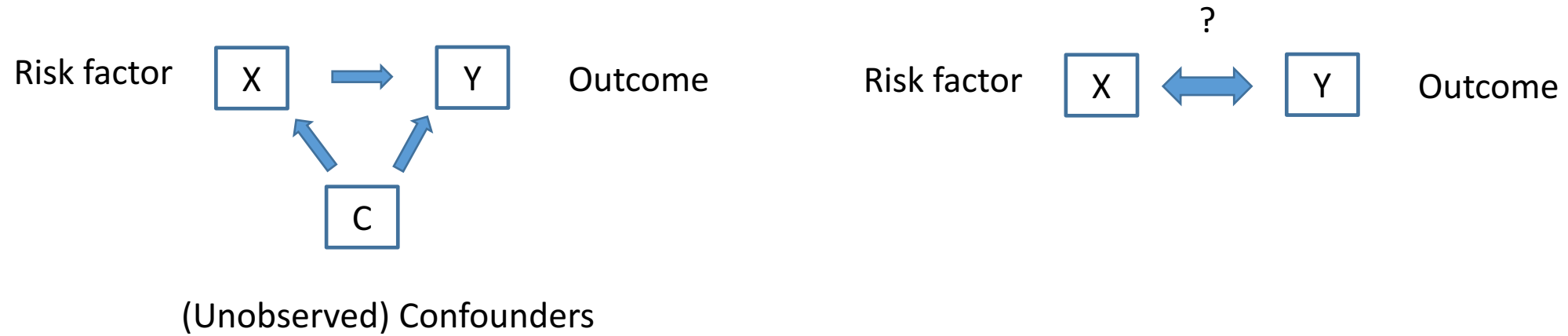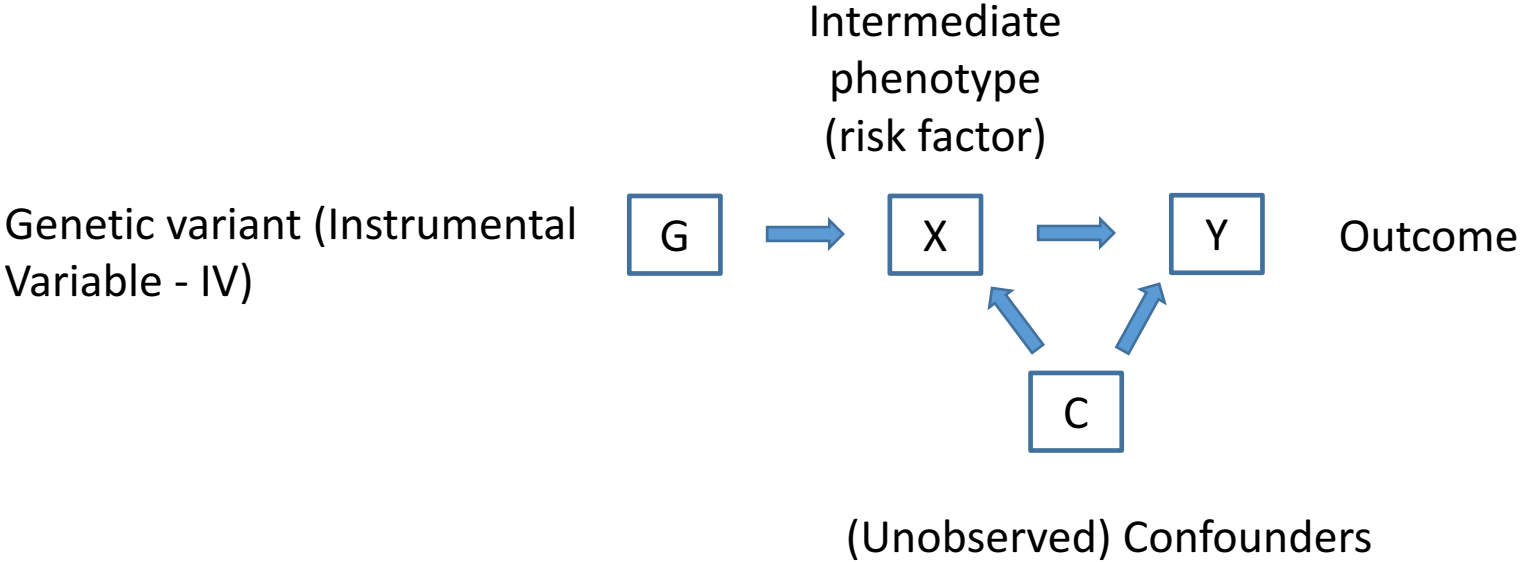# Mendelian Randomization

# Drawback with observational studies
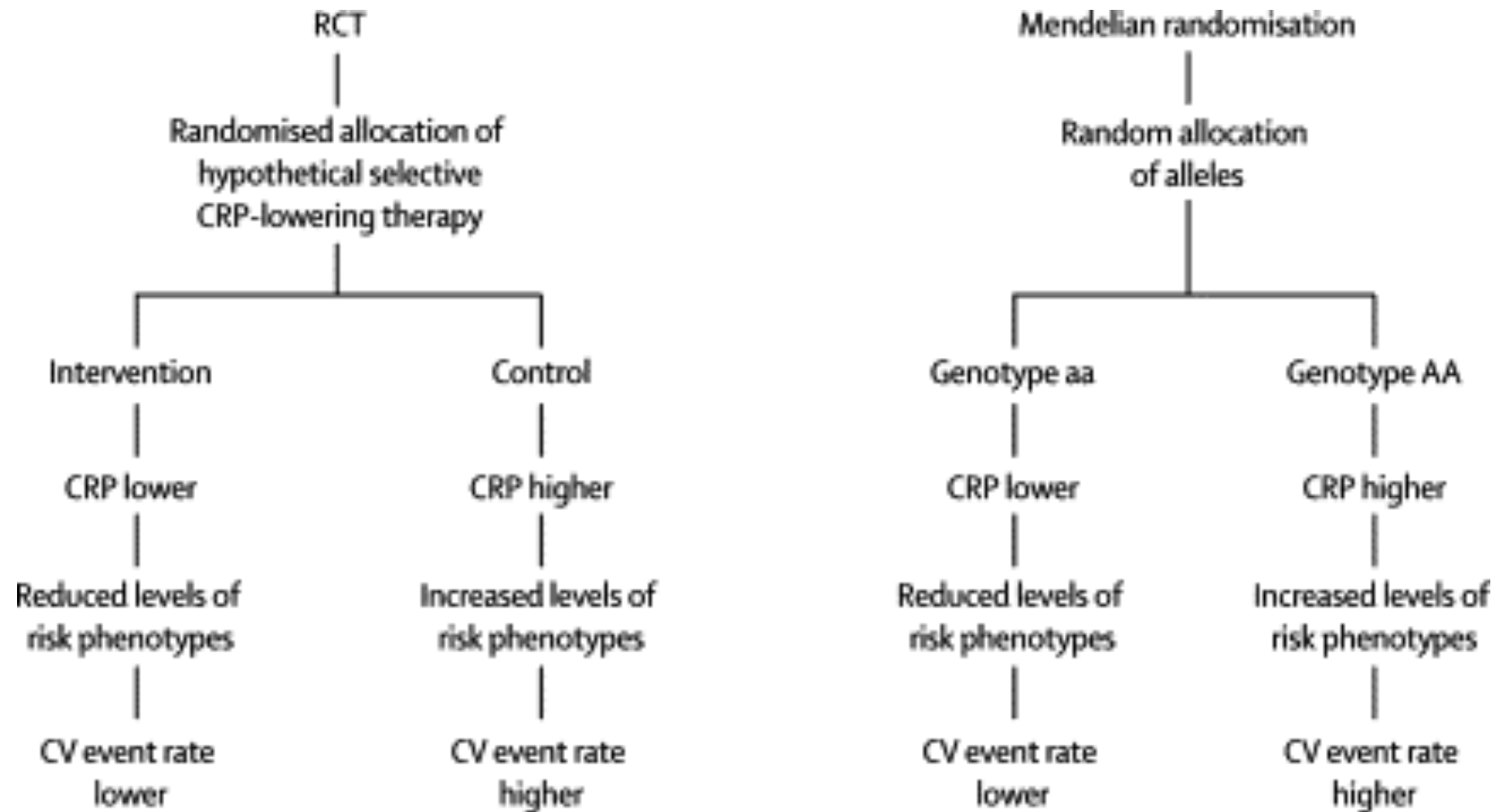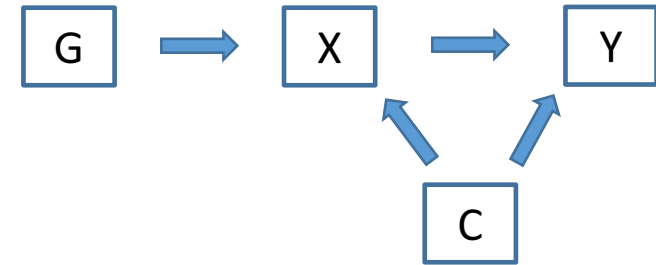
# The power of genetics

# Mendelian Randomization

- *Basic principle: "genetic variants which mirror the biological effects of a modifiable environmental exposure and alters disease risk should be associated with disease risk to the extent predicted by their influence on exposure to the risk factor."*


- The random allocation of genetic variants from parents to offspring means these variants will generally be unrelated to other factors which affect the outcome.


- Furthermore, associations between the genotype and the outcome will not be affected by reverse causation because disease does not affect genotype

Ebrahim & Davey Smith, Hum Genet 2008
Davey Smith & Ebrahim, Int J Epi 2004

Possible effects of C-reactive protein (CRP) on cardiovascular (CV) events. Expected outcome from hypothetical randomized clinical trial of selective CRP-lowering intervention, and from Mendelian randomization analysis, if CRP were causal in developing CV.



Hingorani & Humphries, Lancet 2005

# Three key assumptions in MR analysis

1. G (SNP or a combination of multiple SNPs) is robustly associated with X (risk factor)

2. G is unrelated to any confounders C, that can bias the relationship between G and Y (outcome). In other words, there are no common causes of G and Y (e.g. population stratification)

3. G is related to Y only through its association with X (i.e. no pleiotropy)

# Assumption 1: G is robustly associated with X

- Under certain conditions, the relative bias of the instrument variable (IV) estimate is ~1/F. A "weak" IV has been defined as having F<10, where

$$F = \frac{R^2(n-1-k)}{(1-R^2)k}$$

$R^2$ is variance in X explained by the IV(s), n is sample size and k is number of IVs

- Weak IVs can lead to biased effect estimates (in the direction of the observed X-Y association) in the presence of confounding of the X–Y relationship.
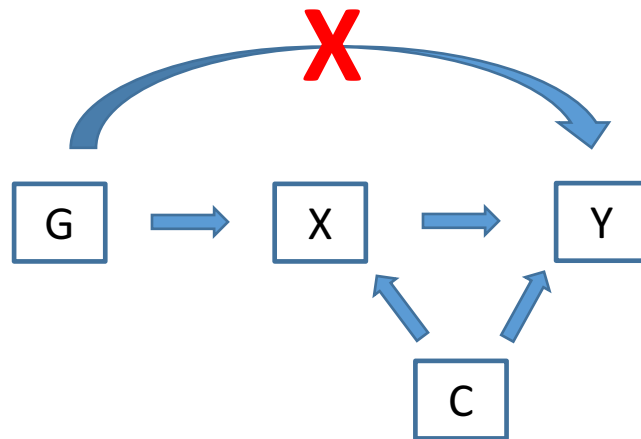
Pierce, IJE 2011

# Assumption 2: No confounding

- G is independent of factors (measured and unmeasured) that confound the X-Y relation

- Since G is randomized at birth and thus is independent of non-genetic confounders and is not modified by the course of disease, the one main concern here is population stratification – i.e. if ancestry is related both to G and Y.

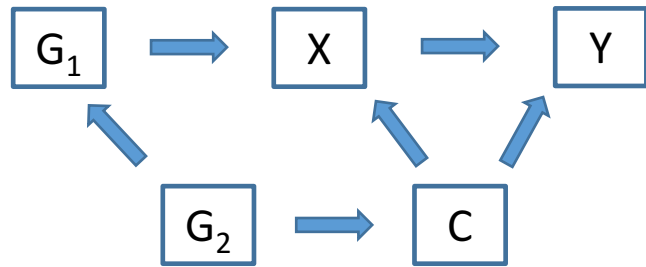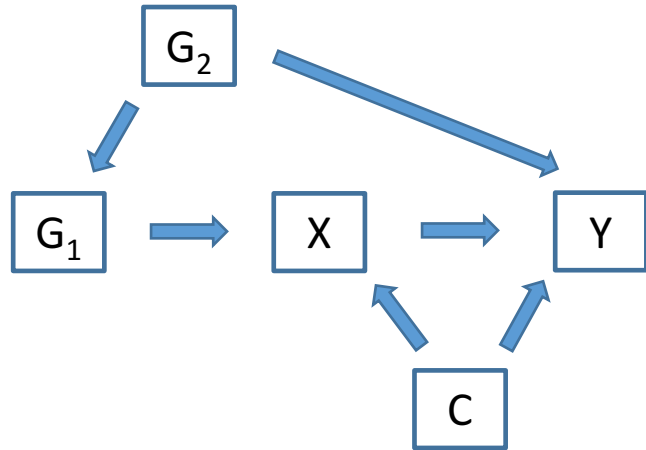- If you have individual-level data, you can test for this (e.g. PCA)

# Assumption 3: No pleiotropy

- This assumption is the trickiest

- Assumes that G is only associated with Y via X and thus the association between G and Y is fully mediated by X and not through any unmeasured factor(s). Needs to be true for SNPs in LD too
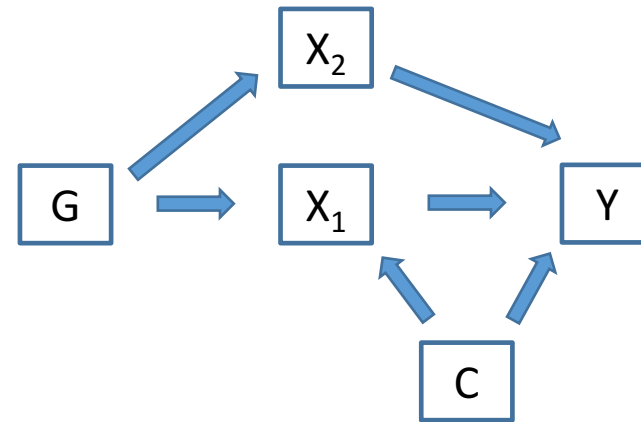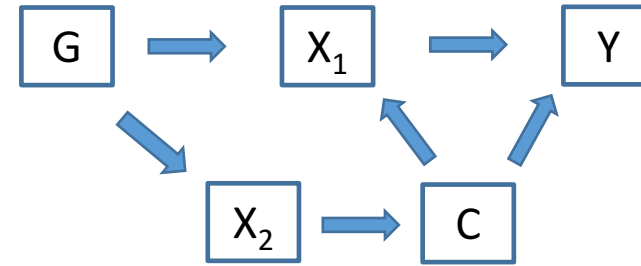
# Scenarios invalidating assumption 3

# Best (but oft-forgotten) practices: the design, analysis, and interpretation of Mendelian randomization studies[1]

*Philip C Haycock,[2]\* Stephen Burgess,[3] Kaitlin H Wade,[2] Jack Bowden,[2,4] Caroline Relton,[2] and George Davey Smith[2]*

**TABLE 2**
Different design strategies for MR[1]

| Study design | Test | Comments |
|---|---|---|
| G-X + G-Y | Implies X→Y | No estimation of magnitude of causal effect |
| One-sample MR | Various hypotheses | Requires individual-level data; lower power; MR estimates are biased toward the confounded observational association by weak instruments |
| Two-sample MR | Various hypotheses | Individual-level or summary data; greater power (due to greater potential sample sizes); MR estimates are biased toward the null by weak instruments |
| Bidirectional MR | X→Y and Y→X | Assesses causation in both directions |
| Two-step MR | X→M→Y | Tests mediation in a causal pathway |
| G×E | X→Y (relation is dependent on environment variable) | Able to detect direct effects (a violation of assumption 2 of MR) |

[1]G×E, gene-environment interaction; G-X, SNP-exposure association; G-Y, SNP-outcome association, M, mediator; MR, Mendelian randomization; SNP, single nucleotide polymorphism; X, hypothesized exposure; Y, outcome variable of interest.

Haycock et al, Am J Clin Nutr 2016

# Individual-level data in one sample

- Access to SNPs, risk factor, and outcome for all participants

- The causal effect of X on Y can be estimated using 2-stage least-squares (2SLS) regression:

1.  X = a + $\gamma G$
2.  Y = c + $\beta X^*$, where $X^*$ are the genetically predicted exposure levels as measured in (1)

- The causal estimate is given by $\beta$
- Can be implemented in R using the "ivpack" package
- Weak instruments cause bias towards the observed confounded association

# Summary data from two samples

- The G-X and the G-Y associations are estimated in two different samples.

- Assumes no overlap among samples and that the two populations are similar (ethnicity, age, sex, etc.)

- Here, bias due to weak IVs will be towards the null

- Note: The G-X and G-Y associations need to be coded using the same effect allele
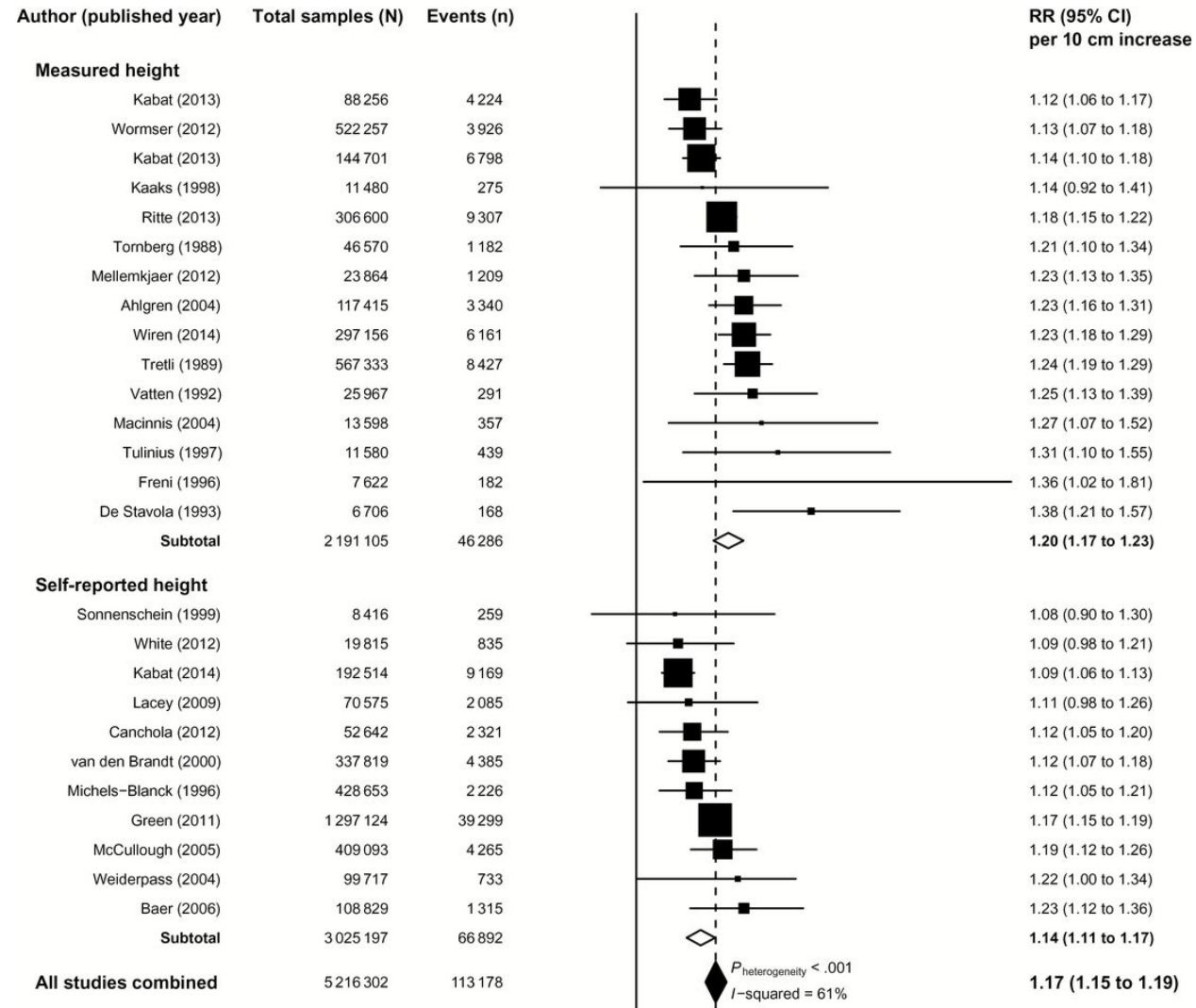
# Summary data from two samples

$$\hat{\beta} = \frac{\sum_{k} \beta_{1k} \beta_{2k} \sigma_{\beta_{2k}}^{-2}}{\sum_{k} \beta_{1k} \sigma_{\beta_{2k}}^{-2}}$$

$$se(\hat{\beta}) = \sqrt{\frac{1}{\sum_{k} \beta_{1k}^{2} \sigma_{\beta_{2k}}^{-2}}}$$

$\beta_{1k}$ is the mean change in X per allele for SNP k, $\beta_{2k}$ is the mean change in Y per allele for SNP k, $\sigma_{2k}^{-2}$ is the inverse variance for the G-Y association.

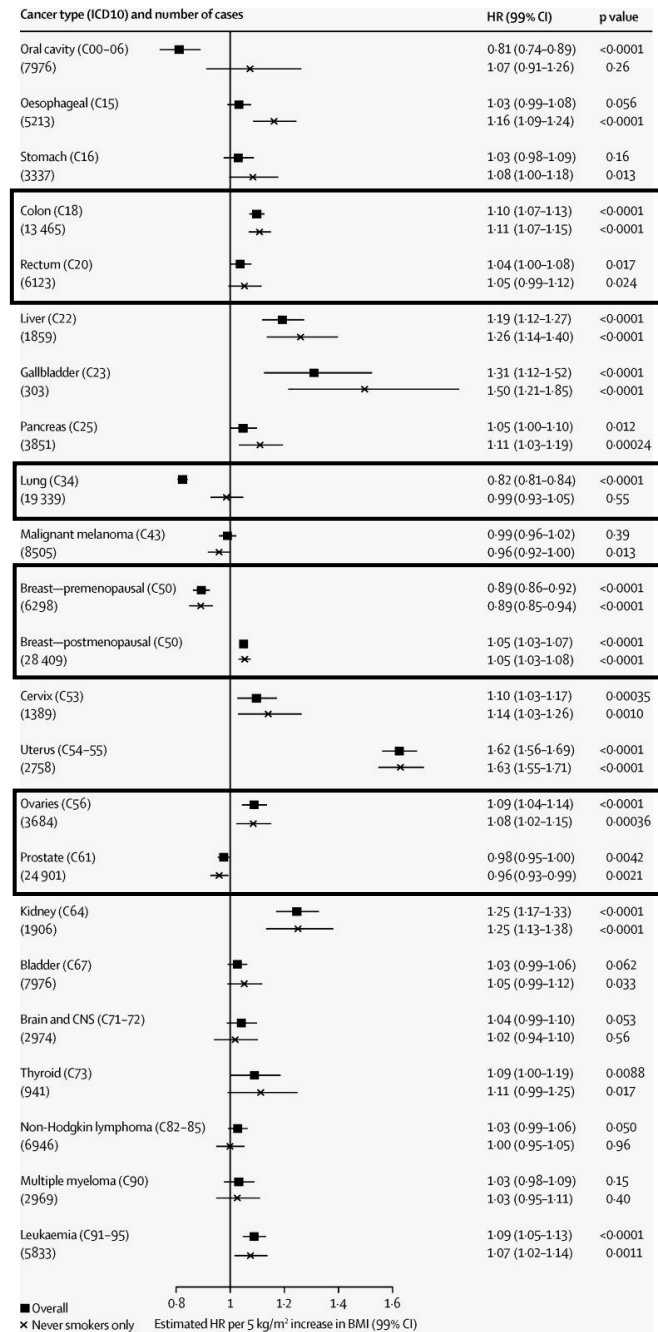# Meta-analysis of associations between height and risk of breast cancer in prospective cohort studies.
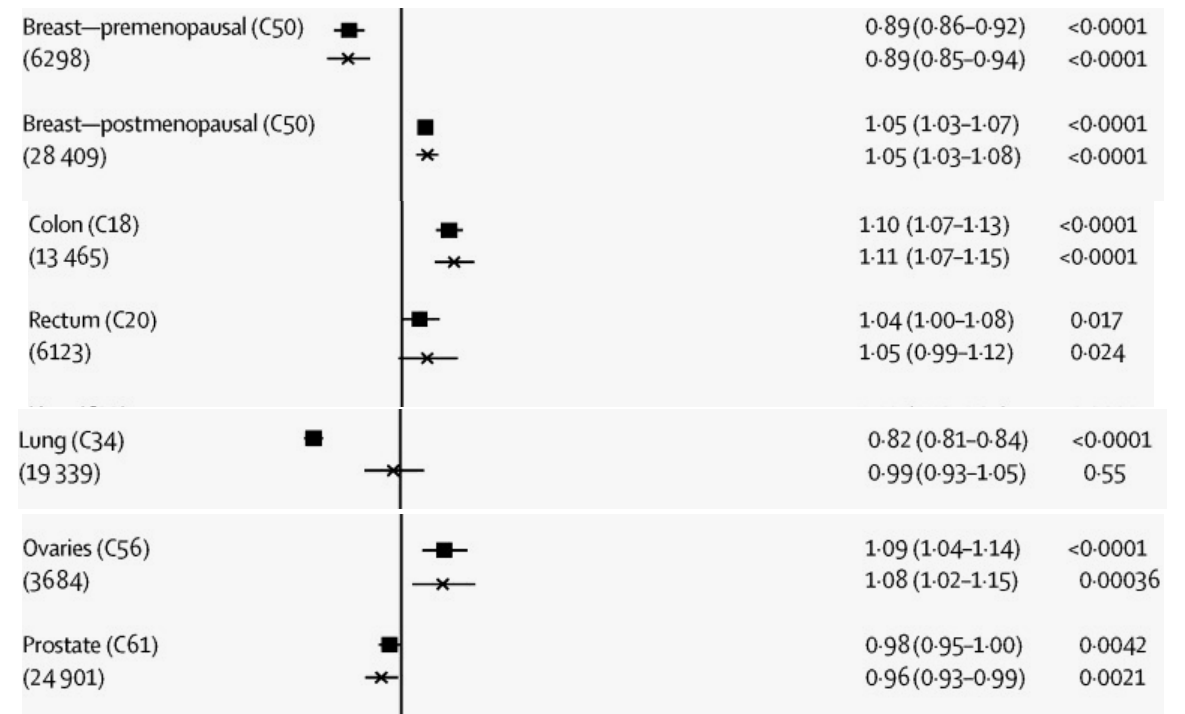
JNCI

**Table 3.**

Association of height and breast cancer risk in women

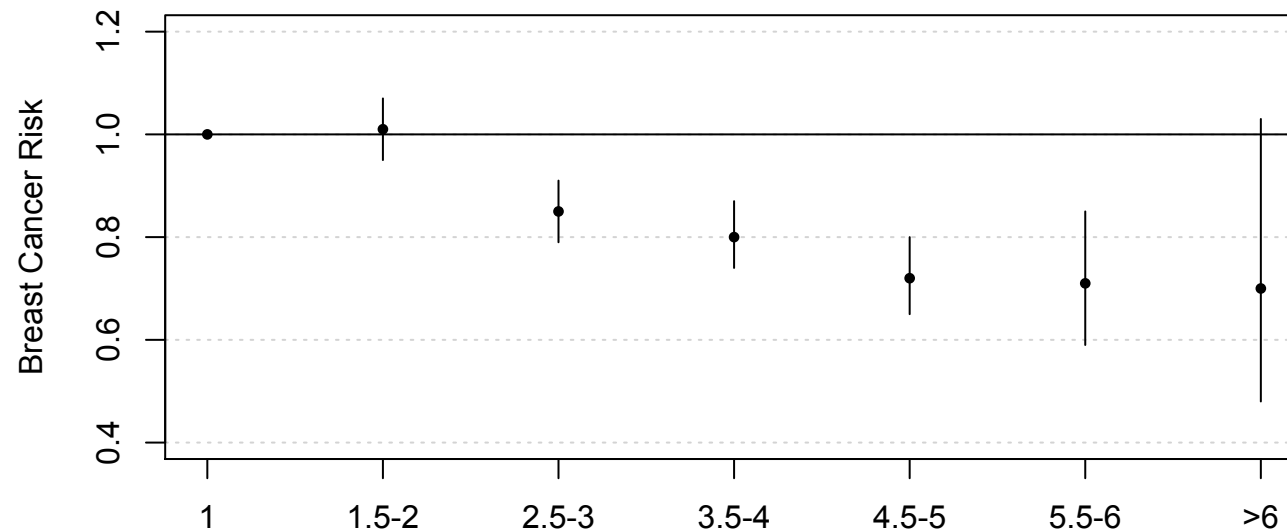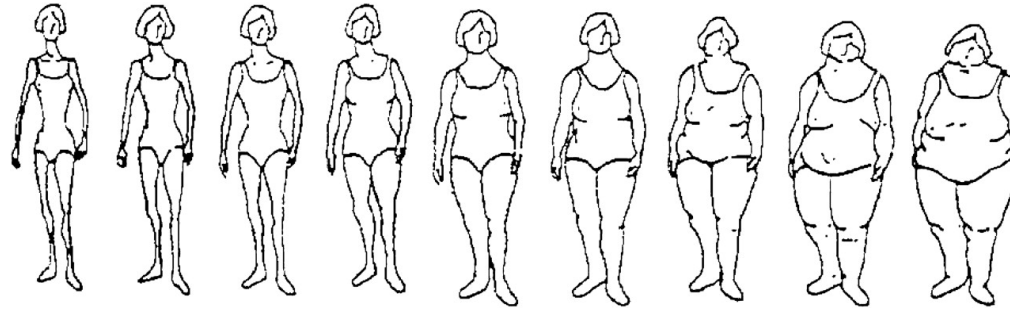| Breast cancer group | Meta-analysis of prospective studies | | | | | | Breast Cancer Association Consortium | | |
| | Observational estimate | | | Instrumental variable estimate | | | Observational estimate | | |
| | N/events | RR (95% CI)* | P | Case patients/ control subjects | OR (95% CI)* | P | Case patients/ control subjects | OR (95% CI)* | P |
|---|---|---|---|---|---|---|---|---|---|
| All women combined | | | | | | | | | |
| All case patients | 5216302/113178 | 1.17 (1.15 to 1.19) | <.001 | 46325/42482 | 1.22 (1.13 to 1.32) | <.001 | 30248/20458 | 1.13 (1.10 to 1.16) | <.001 |
| By menopausal status | | | | | | | | | |
| Premenopausal | 2801907/15439 | 1.16 (1.12 to 1.21) | <.001 | 10209/9053 | 1.29 (1.07 to 1.56) | .007 | 8959/6225 | 1.11 (1.05 to 1.17) | <.001 |
| Postmenopausal | 3111070/63606 | 1.17 (1.14 to 1.21) | <.001 | 23069/19355 | 1.32 (1.17 to 1.49) | <.001 | 20197/13311 | 1.14 (1.10 to 1.18) | <.001 |
| P interaction | | | .79 | | | .86 | | | .35 |
| By ER status | | | | | | | | | |
| ER-positive | 433810/7947 | 1.18 (1.13 to 1.23) | <.001 | 27074/42482 | 1.26 (1.14 to 1.38) | <.001 | 19953/20458 | 1.16 (1.12 to 1.20) | <.001 |
| ER-negative | 433810/1845 | 1.00 (0.87 to 1.13) | .95 | 7288/42482 | 1.02 (0.87 to 1.18) | .84 | 4810/20458 | 1.05 (1.00 to 1.10) | .07 |
| P interaction | | | .02 | | | .02 | | | .002 |

- Association between BMI and cancer risk was assessed for 22 cancers
- 5.24 million individuals (166,996 cancer cases)

Bhaskaran et al, Lancet 2014

# Childhood body fatness is inversely associated with breast cancer risk
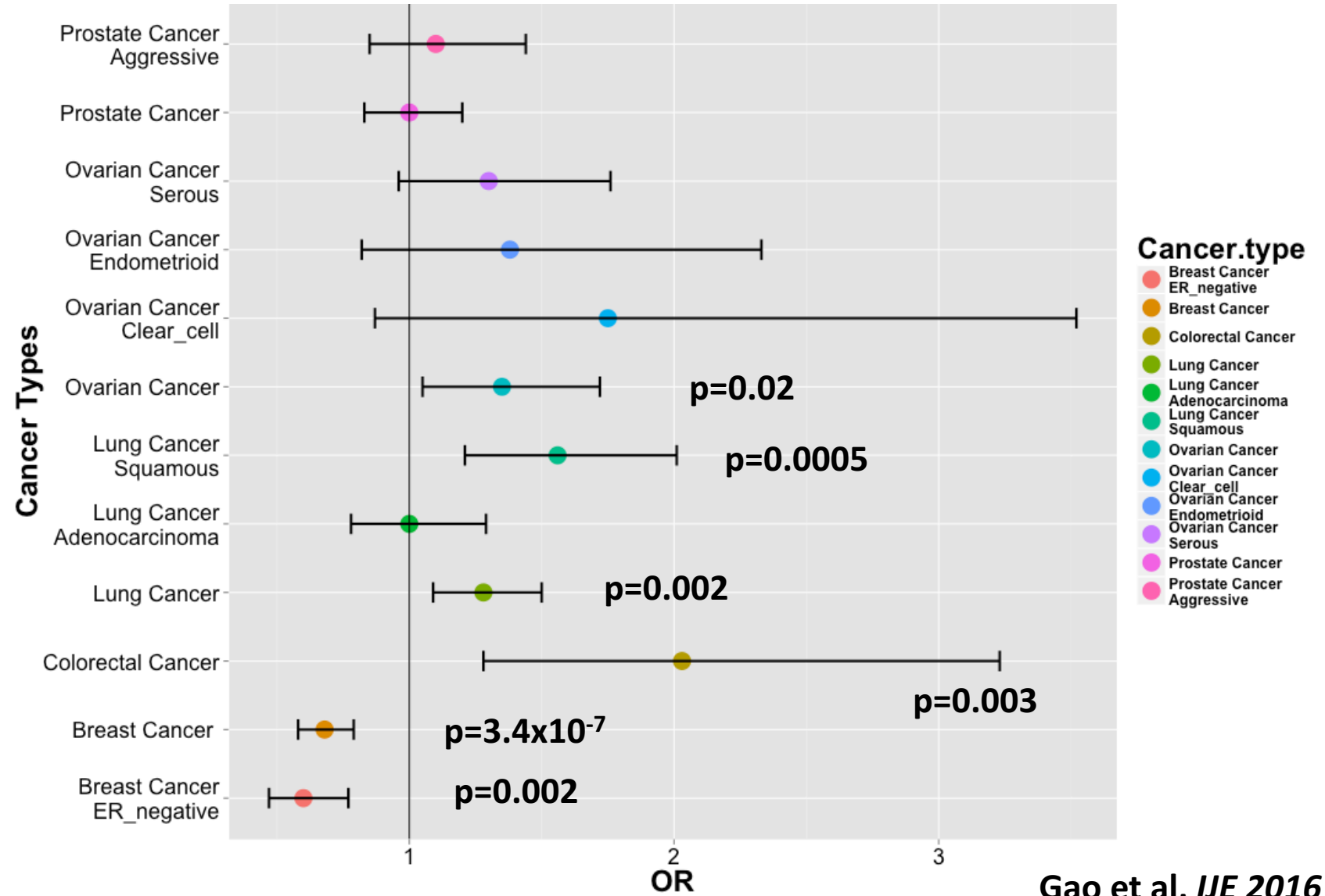
# Expansion to other cancer types within GAME-ON

| Cancer Type | Cases | Controls | GWAS studies |
|---|---|---|---|
| **Breast** | | | |
| All | 15,569 | 18,204 | 11 |
| ER-negative | 4,760 | 13,248 | 8 |
| | | | |
| **Colorectal** | 5,100 | 4,831 | 6 |
| | | | |
| **Lung**[a] | | | |
| All | 12,527 | 17,285 | 6 |
| Adenocarcinoma | 3,804 | 16,289 | 6 |
| Squamous | 3,546 | 16,434 | 6 |
| | | | |
| **Ovarian**[a] | | | |
| All | 4,369 | 9,123 | 3 |
| Clear-cell | 356 | 9,123 | 3 |
| Endometrioid | 715 | 9,123 | 3 |
| Serous | 2,556 | 9,123 | 3 |
| | | | |
| **Prostate** | | | |
| All | 14,160 | 12,712 | 6 |
| Aggressive | 4,446 | 12,724 | 6 |
| **Total** | **51,725** | **62,155** | |

**Gao et al, *IJE 2016***

# Childhood body fatness (9 SNPs)



Gao et al, *IJE 2016*

# Adult BMI (77 SNPs)



Gao et al, *IJE 2016*
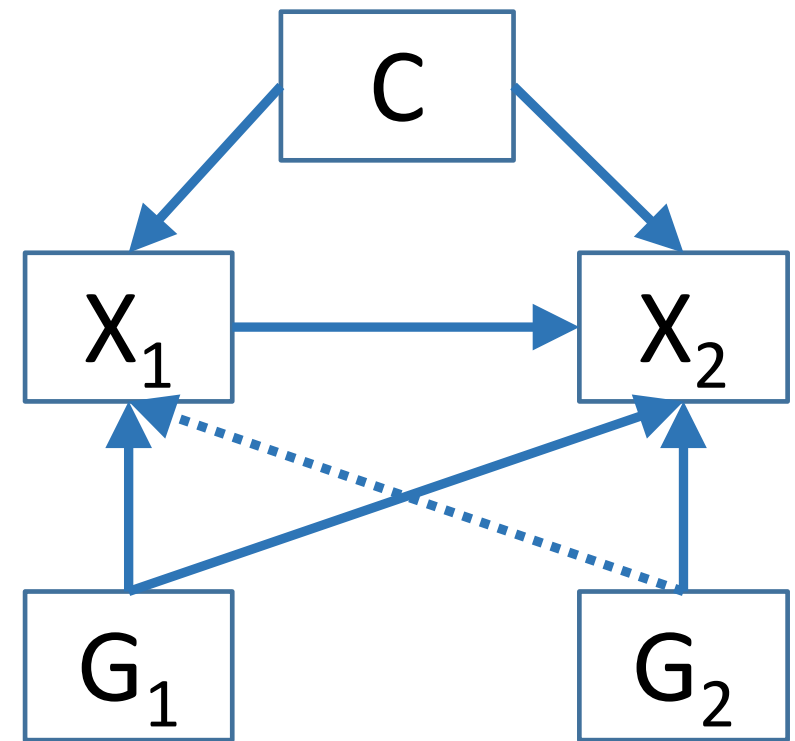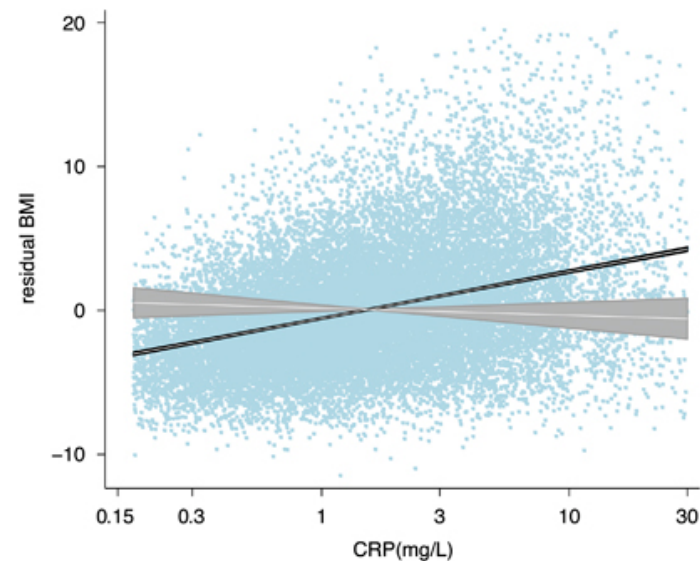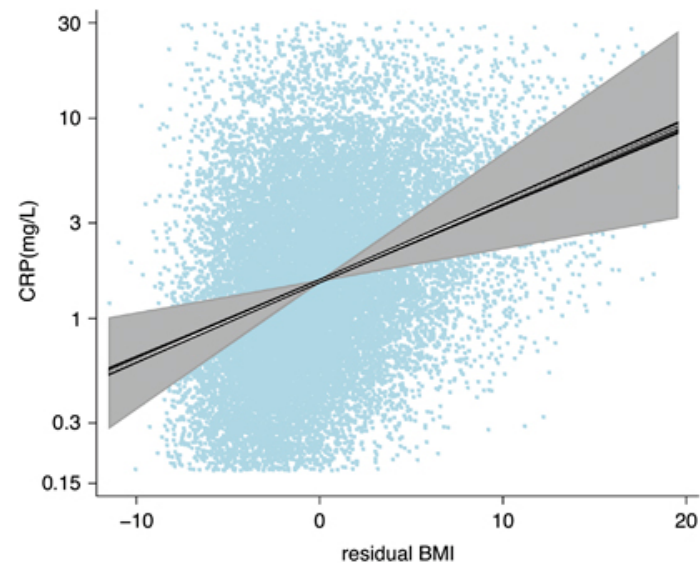
# Bidirectional MR analysis

- Approach to overcome reverse causation

- IVs for both $X_1$ and $X_2$ are used to assess the causal association in both directions

1. Is $G_1$ associated with $X_2$?
2. Is $G_2$ associated with $X_1$?

(Also confirm that $G_1$ is associated with $X_1$ and that $G_2$ is associated with $X_2$

# BMI and CRP – what causes what?

- There is a consistent observed association between high BMI and high CRP levels



Light grey points represent a scatter plot of the correlation between circulating CRP and residual BMI. Gray areas represent 95% confidence regions around IV estimates. Black area represents 95% confidence regions around simple linear regression estimates.

Timpson et al, Int J Obesity 2011

- These data suggest that the observed association between circulating CRP and measured BMI is likely to be driven by BMI, with CRP being a marker of elevated adiposity.

**Table 4. Relationships between genotypic variation and BMI and circulating CRP**

| | | | | *FTO(rs9939609)* | | | | |
|---|---|---|---|---|---|---|---|---|
| | *TT* | *AT* | *AA* | *Per allele effect* | *P-value* | | | |
| BMI | 26.07 (25.98, 26.17) | 26.37 (26.29, 26.45) | 26.73 (26.59, 26.87) | 0.32 (0.24, 0.40) | <0.0001 | | | |
| CRP | 1.51 (1.48, 1.55) | 1.55 (1.52, 1.58) | 1.61 (1.56, 1.67) | 1.03 (1.01, 1.05) | 0.003 | | | |
| | | | | *CRP(rs3091244)* | | | | |
| | *CC* | *CT* | *TT* | *CA* | *AT* | *AA* | *Per allele effect* | *P* |
| BMI | 26.32 (26.23, 26.41) | 26.36 (26.27, 26.44) | 26.24 (26.07, 26.42) | 26.25 (26.02, 26.47) | 26.29 (25.98, 26.61) | 27.15 (26.02, 28.28) | −0.01 (−0.06, 0.04) | 0.7 |
| CRP | 1.37 (1.34, 1.40) | 1.61 (1.57, 1.64) | 1.82 (1.74, 1.90) | 1.71 (1.62, 1.81) | 2.11 (1.95, 2.28) | 2.56 (1.95, 3.37) | 1.11 (1.10, 1.13) | <0.0001 |

**Table 5. Observational and instrumental variable derived relationships between BMI and circulating CRP.**

| Outcome /explanatory variable | Effect estimates | | $P_{IV}$ | $P_{diff}$ | $F_{first}$ |
|---|---|---|---|---|---|
| | Observational | Instrumental variable | | | |
| CRP/BMI | 1.46 (1.44, 1.48) | 1.41 (1.10, 1.80) | 0.006 | 0.8 | 31.1 |
| BMI/CRP | 1.03 (1.00, 1.07) | −0.24 (−0.58, 0.11) | 0.2 | <0.0001 | 57.3 |

Abbreviations: BMI, body mass index; CI, confidence interval; CRP, C-reactive protein.

Observational analysis effects (95% CI) derived from linear regression adjusted for sex, age, age squared, age–sex interaction, log(height), smoking, drinking, education and income.

CRP is log transformed for analyses above and effects on CRP are shown as ratios of geometric means for a s.d. increase in BMI.

BMI effects are expressed as $kg\,m^{-2}$ for a doubling in logCRP.

Instrumental variable derived estimates of the same effects include the same covariates.

$P_{IV}$ is the $P$-value from a test that the instrumental variable estimate is equal to the null.

$P_{diff}$ is the $P$-value from a test for difference between the observational and instrumental variable estimates.

$F_{first}$ is the first stage F-statistic from instrumental variable analysis.

Timpson et al, Int J Obesity 2011

# Drawbacks with MR analysis

- Large sample sizes are needed!

- As genetic effects on risk factors are typically small, MR estimates of association have much wider confidence intervals than conventional epidemiological estimates.

- Make sure that the three key assumptions hold!