

Yale

From cross-phenotype associations to pleiotropy in human genetic studies

Andrew DeWan, PhD, MPH
Associate Professor of Epidemiology
Director, Yale Center for Perinatal, Pediatric and Environmental Epidemiology
Yale School of Public Health

Work done in collaboration with Yasmmy Salinas, PhD, MPH, Assistant Professor of Epidemiology
Yale School of Public Health

Yale SCHOOL OF PUBLIC HEALTH

1

Pleiotropy

- Phenomenon in which a genetic locus affects more than one trait or disease
- Molecular level
 - Single gene with multiple physiological functions
 - Two domains of a single gene product with different functions and affecting multiple phenotypes
 - Gene product with a single function that affects multiple phenotypes acting in multiple tissues
- Statistical level
 - A locus displaying cross-phenotype associations is often considered pleiotropic
 - Can be at the variant, gene or region level

2

Pleiotropy and disease comorbidity

- Examples of correlated (comorbid) disease
 - Obesity, hypertension, dyslipidemia, type 2 diabetes (metabolic disorder)
 - Depression, anxiety, personality disorders (psychiatric disorder)
 - Asthma, obesity (pro-inflammatory conditions)
- Why do certain disease occur together
 - Causality
 - Shared environmental risk factors
 - Shared genetic risk factors

3

Pleiotropy and disease comorbidity

Overlap represents a narrowly-defined phenotype with low heterogeneity (relative to the individual phenotypes)

4

Pleiotropy and disease comorbidity

- Pleiotropy-informed analyses consider multiple phenotypes together and take into account the correlation between the phenotypes
 - Analyzing multiple correlated phenotype (e.g. comorbid diseases) is equivalent to analyzing a single narrowly-defined phenotype with low heterogeneity

5

Pleiotropy and disease comorbidity

- Detecting shared genetics and/or molecular pathways between comorbid diseases can help us understand exactly how the etiology of the diseases overlap
- Etiologic overlaps:
 - provide opportunities for novel interventions that prevent or treat the comorbidity, rather than preventing/treating each disease separately
 - facilitate drug repurposing (that is, known drugs targeting a pleiotropic locus may be repurposed to treat other diseases controlled by that locus, precluding the need for the development and testing of a brand-new drug)

6

Cross-phenotype (CP) associations

Statistical associations between a **single genetic locus** – a single gene or a single variant within a gene – and **multiple phenotypes**

Note that the dashed lines denote uncertainty about whether the SNP has a direct effect on the phenotypes.

13

Analytic options for discovery of CP associations

Univariate **Multivariate**

Key distinction:

- Univariate methods examine the association between a given SNP and each trait *separately*
- Multivariate methods examine the association between a given SNP and each trait by modeling the traits *jointly*

14

Analytic options for discovery of CP associations

Univariate **Multivariate**

Choice between univariate and multivariate approaches depends on:

- Types of data available on our phenotypes of interest
 - Summary statistics vs. individual-level data?
 - Are the phenotypes measured on the same subjects?
- Distribution of the phenotypes (e.g., quantitative or disease trait)

15

Univariate methods are by far the most commonly used to detect CP associations

- Univariate methods include (but are not limited to) the methods you've discussed in class so far:
 - allelic Chi-Square test
 - genotypic Chi-Square test
 - regression-based methods
- The overall approach is to:
 - obtain univariate association p-values for each phenotype
 - declare CP associations at genetic loci that are statistically significantly associated with each phenotype

16

Hypothetical example: Discovery of CP associations for hypertension and heart disease by using logistic regression

Step 1. Fit two univariate regression models within PLINK

$$E[\text{hypertension}] = \beta_0 + \beta_1 * \text{SNP}$$

$$E[\text{heart disease}] = \beta_0 + \beta_1 * \text{SNP}$$

Word of caution: The univariate tests of association should be **marginal tests** (conducted irrespectively of the second phenotype) **NOT conditional tests** (conducted on a subset defined based on absence/presence of the second phenotype). In this example, what that means is that the regression for hypertension should be fit on all subjects *irrespective* of their heart disease status; and the regression for heart disease should be fit on all subjects *irrespective* of their hypertension status. **More on this later!**

evidence to declare a CP association at this SNP.

17

Hypothetical example: Discovery of CP associations for hypertension and heart disease by using logistic regression

Step 1. Fit two univariate regression models within PLINK

$$E[\text{hypertension}] = \beta_0 + \beta_1 * \text{SNP}$$

$$E[\text{heart disease}] = \beta_0 + \beta_1 * \text{SNP}$$

Step 2. For a given SNP, examine p-values for β_1 from each model.

- P-value for β_1 in hypertension model = 1.03×10^{-12}
- P-value for β_1 in heart disease model = 6.02×10^{-9}

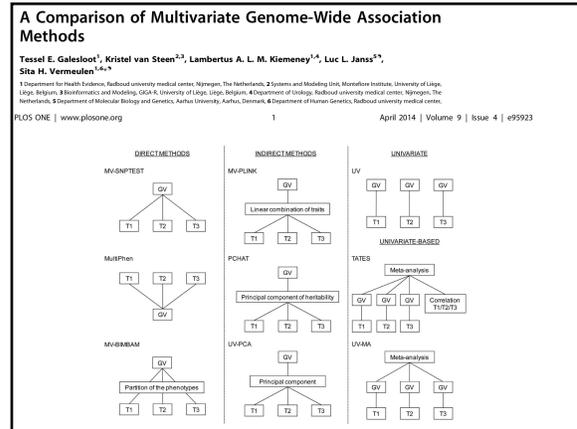
Step 3. Declare CP associations at a given SNP, if the p-values for β_1 in each model surpass the study significance threshold.

- Assuming the standard GWAS significance threshold ($\alpha = 5 \times 10^{-8}$), there is a statistically significant association with both hypertension and heart disease at this particular SNP. Therefore, we have sufficient statistical evidence to declare a CP association at this SNP.

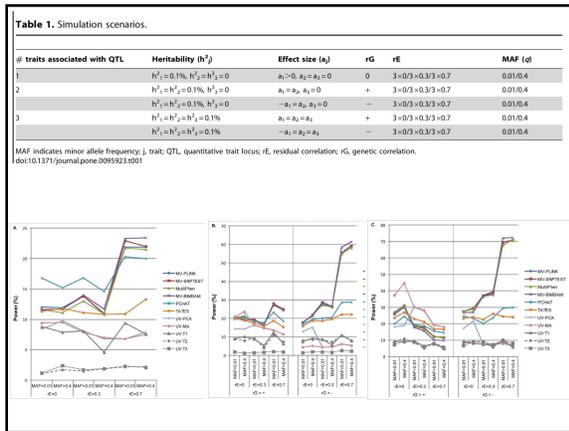
18

Using multivariate methods to increase the power to detect cross-phenotype associations

19



20



21

A comparison of univariate and multivariate GWAS methods for analysis of multiple dichotomous phenotypes

Yasmmy D. Salinas¹, Andrew T. DeWan¹, and Zuoheng Wang²

¹ Department of Chronic Disease Epidemiology, ² Department of Biostatistics, Yale School of Public Health, Yale University, 60 College St, New Haven, Connecticut, USA

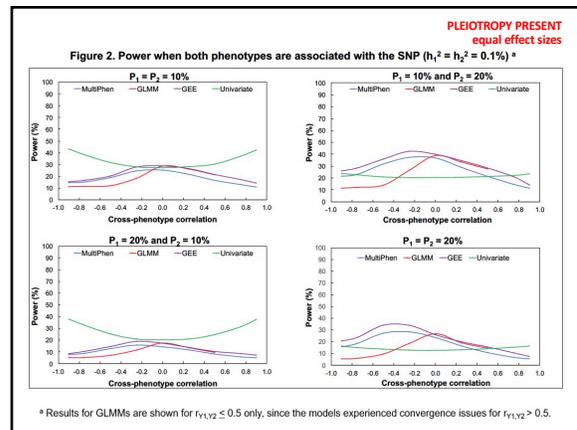
Genet. Epidemiol. 41 (7), 689-689

22

Simulation scenarios

# traits associated	h^2	$r_{Y1,Y2}$	P_j
1	$h_1^2 = 0.1\%$, $h_2^2 = 0\%$	[-0.9, 0.9]	$P_1 = P_2 = 10\%$
			$P_1 = P_2 = 20\%$
			$P_1 = 10\%$, $P_2 = 20\%$
2	$h_1^2 = h_2^2 = 0.1\%$	[-0.9, 0.9]	$P_1 = P_2 = 10\%$
			$P_1 = P_2 = 20\%$
			$P_1 = 10\%$, $P_2 = 20\%$
2	$h_1^2 = 0.1\%$, $h_2^2 = 0.05\%$	[-0.9, 0.9]	$P_1 = P_2 = 10\%$
			$P_1 = P_2 = 20\%$
			$P_1 = 10\%$, $P_2 = 20\%$

23

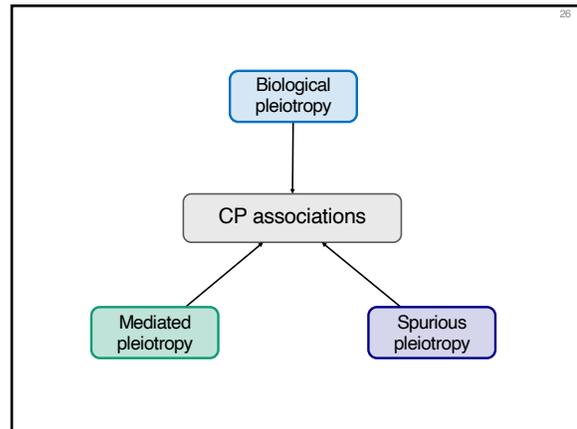


24

25

Problem: CP associations need not be indicative of pleiotropy

25



26

27

Biological pleiotropy

Independent associations between a genetic locus (A) and multiple phenotypic outcomes (Y)

```

    graph LR
      A((A)) --> P1((P1))
      A --> P2((P2))
  
```

The SNP has a direct effect on each phenotype. (Note that direct or causal effects are depicted with solid lines).

27

28

Mediated pleiotropy

Association between a genetic locus (A) and an intermediate phenotype (M) that causes a second phenotypic outcome (Y)

```

    graph LR
      A((A)) --> M((M))
      M --> Y((Y))
  
```

A non-genetic causal link between M and Y induces an association between A and Y, even in the absence of a direct effect of A on Y.

28

29

Spurious pleiotropy

Artificial associations with multiple phenotypes due to issues related to study design, confounding, or associations with markers in strong linkage disequilibrium* with multiple causal variants in different genes

```

    graph LR
      subgraph Left
        C((C)) --> P1((P1))
        C --> P2((P2))
      end
      subgraph Right
        A((A)) --> P1((P1))
        A --> P2((P2))
      end
  
```

*Linkage disequilibrium is the non-random co-segregation of alleles.

29

30

Spurious pleiotropy

Artificial associations with multiple phenotypes due to issues related to study design, confounding, or associations with markers in strong linkage disequilibrium* with multiple causal variants in different genes

```

    graph LR
      subgraph Left
        A((A)) --> P1((P1))
        A --> P2((P2))
        C((C)) --> P1((P1))
        C --> P2((P2))
      end
      subgraph Right
        A((A)) --> P1((P1))
        A --> P2((P2))
        C((C)) --> P1((P1))
        C --> P2((P2))
      end
  
```

Confounders of the relationship between the phenotypes induce spurious cross-phenotype associations

*Linkage disequilibrium is the non-random co-segregation of alleles.

30

Spurious pleiotropy

Artificial associations with multiple phenotypes due to issues related to study design, confounding, or associations with markers in strong linkage disequilibrium* with multiple causal variants in different genes

The SNP has a direct effect on only one of the phenotypes.

*Linkage disequilibrium is the non-random co-segregation of alleles.

31

Spurious pleiotropy

Artificial associations with multiple phenotypes due to issues related to study design, confounding, or associations with markers in strong linkage disequilibrium* with multiple causal variants in different genes

Variables associated with the phenotypes and the SNP induce spurious cross-phenotype associations

*Linkage disequilibrium is the non-random co-segregation of alleles.

32

Spurious pleiotropy

Artificial associations with multiple phenotypes due to issues related to study design, confounding, or associations with markers in strong linkage disequilibrium* with multiple causal variants in different genes

The SNP does not have a direct effect on either phenotype.

*Linkage disequilibrium is the non-random co-segregation of alleles.

33

Guidelines for generating robust statistical evidence of pleiotropy

34

Mediation analysis provides a tool for dissecting CP associations

- Mediation analysis decomposes the **total effect** of the SNP (A) on a phenotypic outcome (Y) into:
 - Direct effect:** effect of A on Y that occurs independently of an intermediate phenotype (M)
 - Indirect effect:** effect of A on Y that occurs through the intermediate phenotype M

35

Mediation analysis: Data requirements

- All phenotypes must be measured on the same subjects
- Temporality must be ascertained
 - The occurrence of the intermediate variable M must precede that of the phenotypic outcome variable Y

36

Mediation analysis: Assumptions

- There must be no unmeasured:
 - confounders of the total effect
 - confounders of the relationship between SNP A and the mediator M
 - confounders of the relationship between mediator M and phenotypic outcome Y

37

Mediation analysis: Assumptions

Typically met in genetic epi studies!

- There must be no unmeasured:
 - confounders of the total effect
 - confounders of the relationship between SNP A and the mediator M
 - confounders of the relationship between mediator M and phenotypic outcome Y

38

Mediation analysis: Assumptions

- There must be no unmeasured:
 - confounders of the total effect
 - confounders of the relationship between SNP A and the mediator M
 - confounders of the relationship between mediator M and phenotypic outcome Y

Requires adjustment for **known** confounders to prevent bias
 (Note: this effectively restricts the use of mediation analyses to datasets in which data on such variables have been collected)

39

Mediation analysis: Regression-based approach

- Requires fitting two regression models, one for mediator M and one for phenotypic outcome Y :
 - $E[M | a, c] = \beta_0 + \beta_1 a + \beta_2' c$
 - $E[Y | a, m, c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_4' c$

Assesses the effect of A on M , while controlling for measured confounders (C)

40

Mediation analysis: Regression-based approach

- Requires fitting two regression models, one for mediator M and one for phenotypic outcome Y :
 - $E[M | a, c] = \beta_0 + \beta_1 a + \beta_2' c$
 - $E[Y | a, m, c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_4' c$

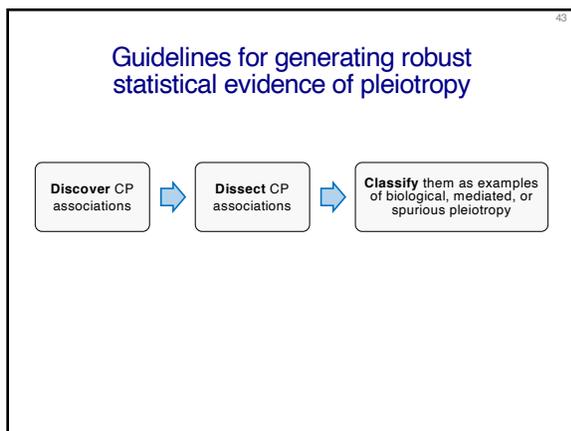
Assesses the effect of A on Y , while controlling for both M and C

41

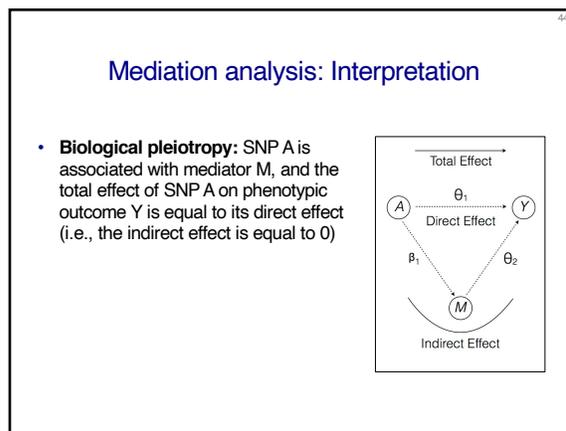
Mediation analysis: Regression-based approach

- Requires fitting two regression models, one for mediator M and one for phenotypic outcome Y :
 - $E[M | a, c] = \beta_0 + \beta_1 a + \beta_2' c$
 - $E[Y | a, m, c] = \theta_0 + \theta_1 a + \theta_2 m + \theta_4' c$
- The parameter estimates from these models (**namely β_1 , θ_1 , and θ_2**) are used to estimate the direct and indirect effects

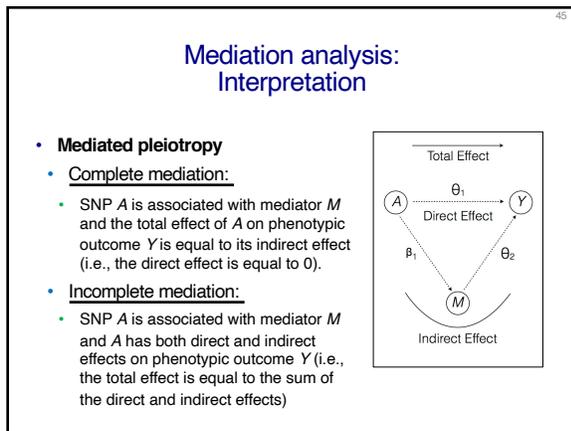
42



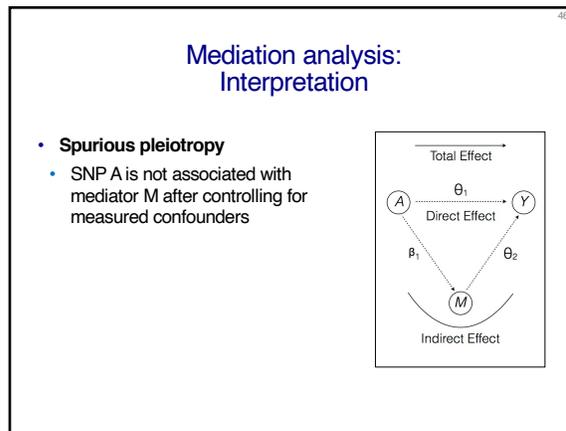
43



44



45



46

mediation R package

```

> med.fit<-glm(W1~rs1_2, data=combined, family=binomial("logit"))
> out.fit<-glm(W2~W1+rs1_2, data=combined, family=binomial("logit"))
> med.out<-mediate(med.fit,out.fit, treat="rs1_2", mediator="W1", boot=TRUE, boot.ci.type="bca", sims=1000)
> summary(med.out)
  
```

Causal Mediation Analysis

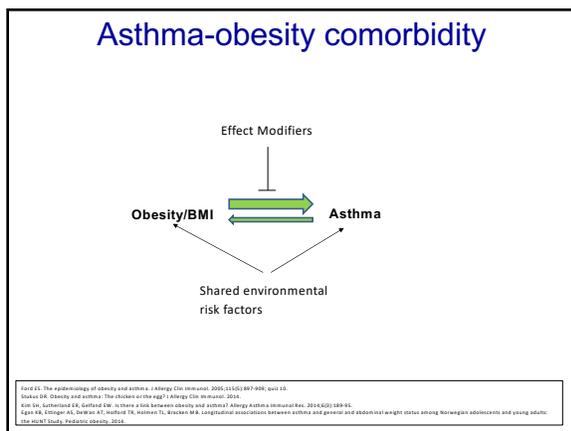
Nonparametric Bootstrap Confidence Intervals with the BCa Method

	Estimate	95% CI Lower	95% CI Upper	p-value
ACME (control)	0.02152	0.01823	0.03	<2e-16 ***
ACME (treated)	0.02199	0.01868	0.03	<2e-16 ***
ADE (control)	0.00723	0.00415	0.01	<2e-16 ***
ADE (treated)	0.00723	0.00443	0.01	<2e-16 ***
Total Effect	-0.00932	-0.02464	-0.09	<2e-16 ***
Prop. Mediated (control)	0.73634	0.65429	0.84	<2e-16 ***
Prop. Mediated (treated)	0.75247	0.67272	0.85	<2e-16 ***
ACME (average)	0.02175	0.01847	0.03	<2e-16 ***
ADE (average)	0.00747	0.00429	0.01	<2e-16 ***
Prop. Mediated (average)	0.74441	0.6654	0.84	<2e-16 ***

47

Empirical searches for pleiotropic loci for asthma and obesity

48



49

Am J Hum Genet. 2009 Jun 8;81(1):87-96. doi: 10.1016/j.ajhg.2009.06.011. Epub 2009 Jul 2.
PRKCA: a positional candidate gene for body mass index and asthma.
 Murthy A¹, Tanfara KJ, Soto-Guerra ME, Avila L, Klanderman BJ, Lake S, Weiss ST, Celedon JC.

Study design

- Two phases:
 - genome-wide linkage analysis of BMI
 - follow-up family-based candidate-gene association study of BMI and asthma
- Strategy for candidate-gene study:
 - Authors focused on a single gene (*PRKCA*) within the BMI linkage peak because:
 - animal models suggest role of *PRKCA* in obesity; and
 - published association studies of other genes within the linkage peak had found no association with BMI.

50

Study population

- Costa Rica study
 - N = 415 asthmatic children + parents
- Childhood Asthma Management Program
 - N = 493 non-Hispanic White asthmatic children + parents

Note that ALL children in both study populations are asthmatic

51

Phenotype definitions

- Body mass index (BMI)
 - calculated from objective measures of height and weight
- Asthma
 - physician-diagnosed asthma + one of the following:
 - 2 respiratory symptoms or asthma attacks in prior year
 - increased airway responsiveness or bronchodilator response

52

Statistical methods

- Univariate family-based association tests (FBATs) were used to test *PRKCA* SNPs for association with BMI and asthma separately
 - Note: The FBAT statistic takes into account the phenotype of the offspring only
- Significance threshold used by study authors: $\alpha = 9.5 \times 10^{-5}$

53

Results for BMI

Table 3. Evidence for Association of *PRKCA* with BMI in Costa Rica and CAMP

Marker	Location (BP) ^a	Minor Allele	Allele Frequency			Number of Informative Families ^b (number of offspring with 0/1 recorded genotype)			Effect Size ^c		CAMP Replication p Value ^{d,e} (two-sided)	Joint p Value ^f (CR, CAMP two-sided)
			CR	CAMP	CR	CAMP	CR	CAMP	CR	CAMP		
rs228853	61874457	T	0.27	0.33	91 (67/24)	110 (80/30)	2.45	1.60	+0.0011	+0.0038 (+0.0076)	5.6×10^{-5} (1.0×10^{-4})	9.5×10^{-5} (1.8×10^{-4})
rs1005651	61868473	C	0.26	0.31	83 (60/23)	113 (83/30)	2.27	1.60	+0.0019	+0.0039 (+0.0077)	9.5×10^{-5} (1.8×10^{-4})	
rs228875	61924337	A	0.29	0.35	101 (70/31)	129 (92/46)	1.71	1.22	+0.0109	+0.0182 (+0.0364)	0.0019 (0.0033)	
rs224497	61931405	C	0.31	0.36	120 (86/34)	136 (98/47)	1.69	1.21	+0.0160	+0.0171 (+0.0341)	0.0025 (0.0046)	

Two BMI-associated variants

54

Results for asthma

Table 4. Evidence for Association of *PRKCA* with Asthma in Costa Rica and CAMP

Marker	Location (BP) ^a	Minor Allele	Allele Frequency		Number of Informative Families ^b (number of offspring with 0/1 recorded genotype)		Costa Rica p Value ^{c,d}	CAMP Replication p Value ^{c,d} (two sided)	Joint p Value ^e (CR, CAMP two-sided)
			CR	CAMP	CR	CAMP			
rs732191	61779673	G	0.46	0.35	168 (117/51)	141 113/43	-0.0194	-0.0214 (-0.0428)	0.0036 (0.0067)
rs9895580	61789701	C	0.47	0.35	168 (117/51)	141 114/43	-0.0171	-0.0160 (-0.0320)	0.0025 (0.0047)
rs4411531	61793662	A	0.29	0.12	88 (70/18)	25 (24/1)	-0.0058	-0.0058 (-0.0117)	0.0004 (0.0007)
rs8080771	61824330	G	0.46	0.35	164 (116/48)	108 (90/29)	-0.0161	-0.0070 (-0.0140)	0.0011 (0.0021)
rs11652956	61839798	G	0.29	0.12	83 (65/18)	23 (22/1)	-0.0101	-0.0111 (-0.0222)	0.0011 (0.0021)
rs7221968	61848731	C	0.27	0.11	79 (63/16)	18 (17/1)	-0.0122	-0.0216 (-0.0432)	0.0024 (0.0045)
rs7405806	61862056	A	0.49	0.31	164 (109/55)	90 (77/20)	-0.0309	-0.0009 (-0.0018)	0.0003 (0.0006)
rs11079657	61862528	A	0.38	0.23	129 (94/35)	60 (56/8)	-0.0092	-0.0002 (-0.0004)	2.6 × 10⁻¹⁰** (3.0 × 10 ⁻¹⁰ ***)

One asthma-associated variant

Conclusions

- Authors' conclusion: *PRKCA* displays pleiotropy for asthma and BMI (pleiotropy at gene level)
- Two variants (rs228883 and rs1005651) displayed statistically significant associations with body mass index
- A different variant (rs11079657) displayed a statistically significant association with asthma.

Conclusions

- Our conclusion: *PRKCA* is associated with asthma and with BMI among asthmatics (no true CP association!)
- There is insufficient evidence to declare a CP association at *PRKCA* because the test of association with BMI was not a marginal test
 - FBAT test for BMI only took into account the phenotype of the offspring – which were ALL asthmatic
- Thus, it remains to be seen whether the association with BMI is also present among non-asthmatics subjects
- Without that information, we would not be able to assess whether asthma is a **mediator** or a **moderator** of the relationship between *PRKCA* and BMI.

A GWAS study of pleiotropy

Discovery and Mediation Analysis of Cross-Phenotype Associations Between Asthma and Body Mass Index in 12q13.2

Yasmmyn D. Salinas¹, Zuohe Wang, and Andrew T. DeWan

¹ Correspondence to Dr. Yasmmyn D. Salinas, Department of Chronic Disease Epidemiology, Yale School of Public Health, 60 College Street, New Haven, CT 06520 (e-mail: yasmmyn.salinas@yale.edu).

Am J Epidemiol. 2021;190(1):85-94

Study design

- Two parts:
 - Genome-wide search for cross-phenotype associations with asthma and body mass index
 - Follow-up mediation analysis to dissect genome-wide significant CP associations

Study population

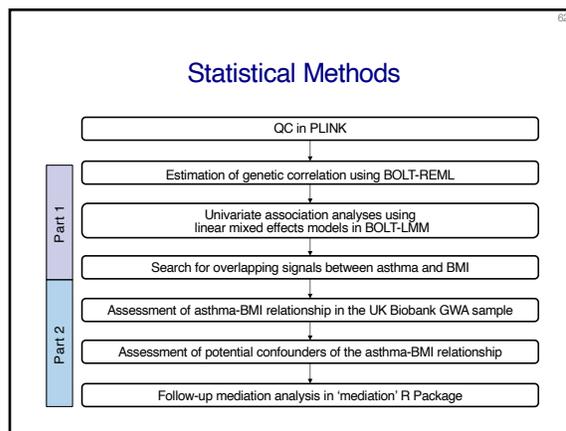
- N = 305,945 White, British subjects from the UK Biobank (a population-based prospective cohort study of > 500,000 subjects, aged 40-69 years at baseline)

Phenotype definitions

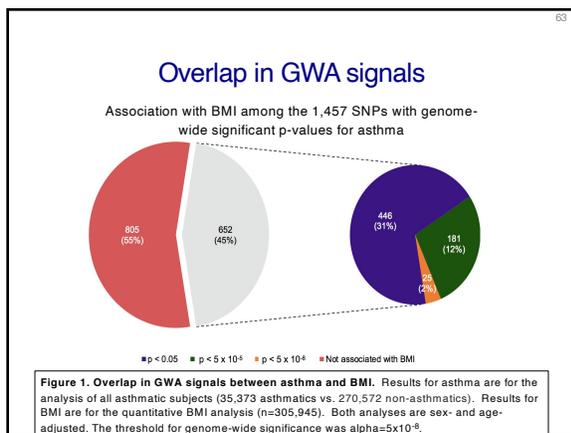
- BMI at baseline (kg/m²):
 - calculated based on height and weight measurements collected by trained UK Biobank staff at the recruitment sites
- Asthma diagnosed prior to baseline (yes/no):
 - ascertained via the question "Has a doctor ever told you that you had asthma?"
 - Note: In mediation analyses, two subgroups were created based on age-at-diagnosis



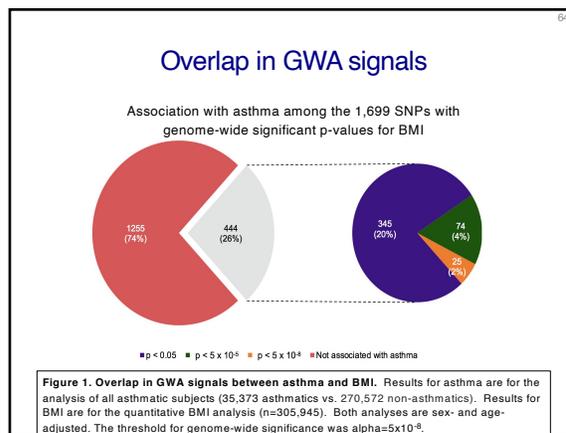
61



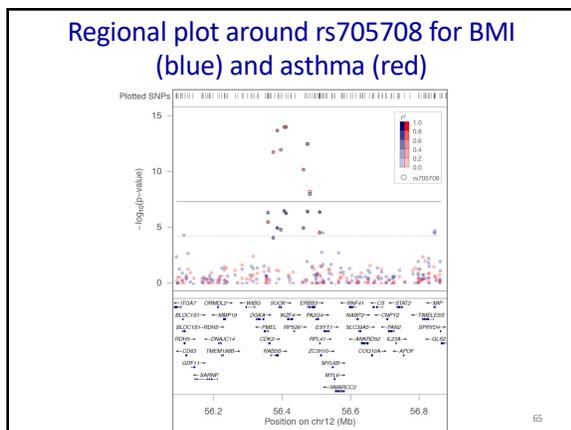
62



63



64



65

Cross-phenotype associations in 12q13.2

Table 2. Cross-phenotype associations in 12q13.2^a

SNP	Gene	BP	Effect/reference allele	EAF	OR (95% CI)	P ^b	BMI ^c	P ^d
rs2069408	<i>CD42</i>	56,364,321	G/A	0.3388	1.04 (1.02, 1.06)	3.30x10 ⁻⁵	-0.06 (-0.08, -0.04)	5.40x10 ⁻⁷
rs1873914	<i>RAB3</i>	56,379,427	C/G	0.4237	1.06 (1.04, 1.08)	2.40x10 ⁻⁵	-0.05 (-0.07, -0.02)	7.90x10 ⁻⁷
rs705708	<i>SUGX</i>	56,390,656	G/A	0.3376	1.07 (1.05, 1.09)	3.10x10 ⁻⁵	-0.05 (-0.08, -0.03)	1.10x10 ⁻⁶
rs10876864 ^a	<i>SUGX</i>	56,401,085	G/A	0.4279	1.06 (1.04, 1.08)	1.50x10 ⁻⁵	-0.05 (-0.07, -0.03)	1.60x10 ⁻⁶
rs1701704	<i>RZFP4</i>	56,412,487	G/T	0.3433	1.07 (1.05, 1.09)	1.50x10 ⁻⁵	-0.06 (-0.09, -0.04)	3.70x10 ⁻⁶
rs2456973	<i>RZFP4</i>	56,416,928	C/A	0.3432	1.07 (1.05, 1.09)	1.50x10 ⁻⁵	-0.06 (-0.08, -0.04)	6.00x10 ⁻⁶
rs11171739 ^a	<i>ERBB3</i>	56,470,625	C/T	0.4337	1.06 (1.04, 1.07)	8.80x10 ⁻⁶	-0.05 (-0.07, -0.03)	1.10x10 ⁻⁶
rs2292239	<i>ERBB3</i>	56,482,180	T/G	0.3470	1.07 (1.05, 1.08)	4.50x10 ⁻⁶	-0.06 (-0.08, -0.04)	4.20x10 ⁻⁶
rs705708	<i>ERBB3</i>	56,488,913	A/G	0.4712	1.05 (1.03, 1.07)	7.20x10 ⁻⁶	-0.06 (-0.09, -0.04)	1.30x10 ⁻⁶
rs11171743 ^a	<i>ESTY7</i>	56,518,408	T/G	0.6180	1.04 (1.02, 1.05)	2.90x10 ⁻⁶	-0.06 (-0.08, -0.04)	4.50x10 ⁻⁶

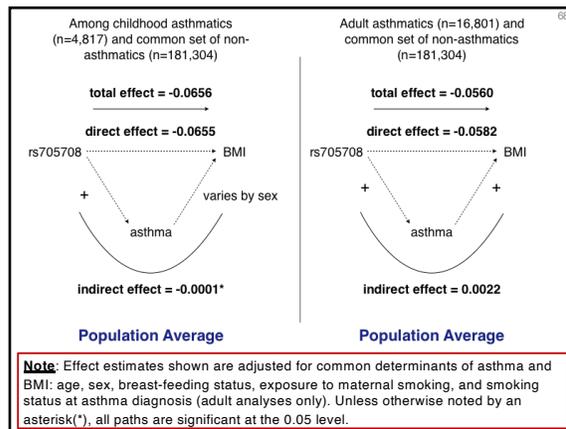
Abbreviations: BP = base-pair; BMI = body mass index; CI = confidence interval; EAF = effect allele frequency; OR = odds ratio; SNP = single-nucleotide polymorphism

^a Results shown for SNPs with p < 5x10⁻⁶ for asthma and p < 0.05 for BMI.
^b For intergenic SNPs, the nearest gene is listed, with priority given to genes directly downstream of variant.
^c P-value from BOLT-LMM, derived using the standard "infinite" mixed model.
^d P-value from BOLT-LMM, derived using the Gaussian mixture model.

66

Decomposing the effect of rs705708 on BMI via mediation analysis

67



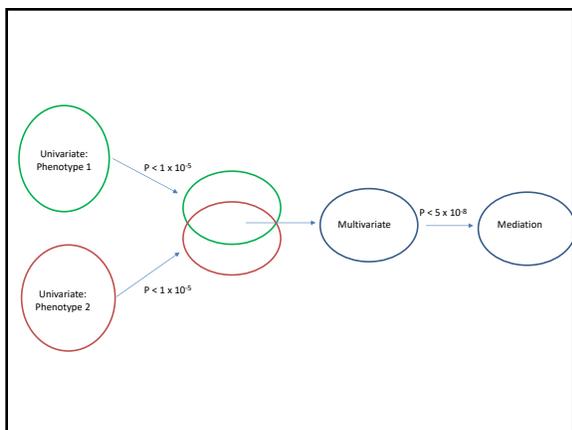
68

- ### Conclusions
- rs705708 has a positive direct effect on asthma
 - Stronger in magnitude for childhood asthma
 - rs705708 has a negative direct effect on BMI
 - Consistent in magnitude and direction in analyses including childhood vs. adult asthmatics
 - This suggests that locus 12q13.2, tagged by rs705708, has pleiotropic effects on asthma and BMI.

69

- ### Conclusions
- 12q13.2 is multigenic and our CP associations span genes *CDK2*, *RAB5*, *SUOX*, *IZK4*, *RPS26*, *ERBB3*, and *ESYT1*.
 - rs705708 is the top regional BMI signal and resides in *ERBB3*.
 - The top regional asthma signal, rs2456973, resides in *IZKF4*.
 - While rs705708 and rs2456973 could be in LD with the same causative variant in either *ERBB3* or *IKZF4* or another gene in 12q13.2, it is also possible that each variant could tag a distinct, trait-specific causative variant in different genes.
 - Therefore, locus 12q13.2 displays pleiotropic effects on asthma and BMI, but this may not be an example of pleiotropy at the gene level (biological pleiotropy).

70



71