

1 **Mega-analysis quantifies ‘hidden heritability’ due to heterogeneity**

2 **ACCEPTED FOR PUBLICATION NATURE HUMAN BEHAVIOUR**

3 Felix C. Tropf^{1*}, S. Hong Lee², Renske M. Verweij³, Gert Stulp³, Peter J. van der Most⁴,
4 Ronald de Vlaming^{5,6}, Andrew Bakshi⁷, Daniel A. Briley⁸, Charles Rahal¹, Robert Hellpap¹,
5 Anastasia Nyman⁹, Tõnu Esko^{10,11}, Andres Metspalu¹⁰, Sarah E. Medland¹², Nicholas G.
6 Martin¹², Nicola Barban¹, Harold Snieder⁴, Matthew R. Robinson^{7,13}, Melinda C. Mills¹

7
8 ¹Department of Sociology/ Nuffield College, University of Oxford, Oxford OX1 3UQ, UK

9 ²School of Environmental and Rural Science, The University of New England, Armidale
10 NSW 2351, Australia

11 ³Department of Sociology/Interuniversity Center for Social Science Theory and
12 Methodology, University of Groningen, Groningen 9712 TG, The Netherlands

13 ⁴Department of Epidemiology, University of Groningen, University Medical Center
14 Groningen, Groningen 9700 RB, Netherlands.

15 ⁵Erasmus University Rotterdam Institute for Behavior and Biology, Erasmus School of
16 Economics, Rotterdam, the Netherlands

17 ⁶Department of Applied Economics, Erasmus School of Economics, Rotterdam, the
18 Netherlands

19 ⁷Institute of Molecular Biosciences, The University of Queensland, Brisbane, QLD
20 4072, Australia

21 ⁸Department of Psychology, University of Illinois at Urbana-Champaign, Champaign 61820-
22 9998, USA

23 ⁹Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, PO Box 281,
24 Stockholm SE-171 77, Sweden

25 ¹⁰Estonian Genome Center, University of Tartu, Tartu, Estonia, 51010

26 ¹¹Broad Institute of the Massachusetts Institute of Technology and Harvard
27 University, 140 Cambridge 02142, MA, USA

28 ¹²Quantitative Genetics Laboratory, QIMR Berghofer Medical Research Institute, Brisbane,
29 QLD 4029, Australia

30 ¹³Department of Computational Biology, University of Lausanne, Lausanne, CH-1015,
31 Switzerland

32

33

34 Corresponding author: Felix C. Tropf, Manor Rd, Oxford, UK, OX1 3UQ, +44 (0) 1865 28
35 17 40, felix.tropf@sociology.ox.ac.uk

36

37

1 **Abstract**

2 The genome-wide association studies (GWAS) dominating genetic discovery are based on
3 large meta-analyses that combine data from diverse historical time periods and populations.
4 Polygenic scores constructed from GWAS explain only a fraction of the heritability derived
5 from twin and whole-genome studies on single-populations, known as the ‘hidden heritability’
6 puzzle. Using seven sampling populations (N=35,062), we test whether hidden heritability is
7 attributed to genetic heterogeneity, showing that estimates are substantially and generally
8 smaller from across compared to within populations. We show that the hidden SNP-based
9 heritability ranges from zero (height), 20% (BMI), 37% (education), 40% (age at first birth) to
10 75% (number of children). Simulations demonstrate that our results more likely reflect
11 heterogeneity in phenotypic measurement or gene-environment interaction than genetic
12 explanations. These findings have substantial implications for gene discovery, suggesting that
13 large homogenous datasets are required for behavioural phenotypes and that gene-environment
14 interaction is a central challenge for genetic discovery.

15

16 Keywords: human reproduction; age at first birth; educational attainment; gene-environment
17 interaction; missing heritability; hidden heritability

18

1

2 **Introduction**

3 The genome-wide association studies (GWAS) that dominate genetic discovery are based on
4 diverse samples that span vast historical time periods and populations.¹ The proportion of
5 phenotypic variance accounted for by single-nucleotide polymorphisms (SNPs) that reach
6 genome-wide significance, and the polygenic scores constructed from all SNPs using GWA
7 study results, however, represent only a fraction of heritability estimates derived from twin and
8 other whole-genome studies.²⁻⁴

9 To understand this disparity, it is essential to explain three central ways to measure
10 heritability (see Box 1 for detailed definitions). First, narrow-sense heritability stems from
11 family-based studies and often twin research (h^2_{family}) and produces the highest heritability
12 estimates. These studies demonstrated a genetic basis for anthropometric traits such as height
13 and body mass index (BMI), but also behavioral phenotypes such as educational attainment
14 and human reproductive behavior (i.e., number of children, age at first birth).⁵⁻⁷ A recent meta-
15 analysis of twin studies from 1958-2012⁵ estimated, for instance, heritability for educational
16 attainment as 52% (N=24,484 twin pairs) and 31% for reproductive traits (N=28,819 twin
17 pairs).

18 GWAS heritability estimates (h^2_{GWAS}) use whole genome-data to estimate the proportion
19 of phenotypic variance accounted for by genetic variants known to be robustly associated with
20 the phenotype of interest and produce the lowest estimates. The polygenic score from a recent
21 meta-GWAS of educational attainment with over 300,000 participants, explains around 4% of
22 the variance⁴ with another GWAS for age at first birth explaining only 1%.²

1 Yang and colleagues argued that most genetic effects are too small to be reliably
2 detected in GWAS of current sample sizes and proposed an alternative approach: whole
3 genome restricted maximum likelihood estimation (GREML) performed by GCTA software.^{8,9}
4 This third measure is often referred to as SNP- or chip-based heritability (denoted by h^2_{SNP}), and
5 is the proportion of phenotypic variance explained by additive genetic variance jointly
6 estimated from all common variants on standard GWAS chips. These estimates are typically
7 between h^2_{family} and (h^2_{GWAS}) estimates. Contrary to the low h^2_{GWAS} estimates of between 1–4% for
8 these phenotypes, the SNP-heritability for educational attainment has been estimated as 22%,
9 15% for age at first birth.^{3,10}

10 This stark discrepancy in heritability estimates has spawned debates about ‘missing
11 heritability’ (difference between h^2_{GWAS} and h^2_{SNP} with h^2_{family}) and ‘hidden heritability’
12 (difference between whole-genome derived estimates from h^2_{GWAS} and h^2_{SNP}) (for full
13 definitions see Box 1 and¹¹).^{12–15} ‘Missing heritability’ has been linked to fundamental
14 differences in study designs between family and whole-genome studies¹² and that non-additive
15 genetic effects^{13,14} and inflated estimate from twin studies due to shared environmental factors
16 ¹⁶ might explain part of this phenomenon. , Empirical studies find no evidence for either of
17 these reasons. A recent investigation on height and BMI, however, demonstrates that the
18 inclusion of rare genetic variants can strongly influence the heritability estimate.¹⁵ The
19 underlying reason for the discrepancy of ‘hidden heritability’ between h^2_{SNP} versus h^2_{GWAS}
20 estimates, are less well understood.¹⁷

21 Here, we interrogate the common assumption underlying GWA studies’ meta-analyses,
22 that genetic effects are ‘universal’ across environments. The large GWAS meta-analyses
23 required to detect SNP associations consist of a wide array of samples across historical periods

1 and countries, representing heterogeneous populations subject to diverse environmental
2 influences. Heterogeneity across environments can emerge for different reasons such as
3 differences in population structure, genotype or phenotype measurement, heterogeneous
4 imputation quality across populations or sensitivity of the phenotype to environmental change.
5 Demographic research has shown, that education and reproductive behavior is strongly
6 modified by environmental changes such as female educational expansion or the introduction
7 of effective contraception.¹⁸ If genetic effects are not universal but rather heterogeneous across
8 populations, heritability estimates from GWAS meta-analyses should produce weaker signals
9 and we would witness a reduction in both the discovery rate and the variance explained from
10 SNPs across populations.¹⁹

11 We conduct a mega-analysis using whole-genome methods, which entails pooling all
12 cohorts to estimate genetic relatedness not only within, but also across populations. We first
13 apply models based on GREML estimation⁹ on primary data from seven pooled sampling
14 populations, which estimates the average common SNP-based heritability (h^2_{SNP}) between and
15 within environments. We subsequently apply a gene-environment interaction models adding a
16 within population matrix to estimate the average SNP-based heritability within populations in
17 our data and decompose the variance explanation of common SNPs within and between
18 sampling populations and birth cohorts.^{9,20} If SNP-based heritability is significantly higher
19 within than across environments, we conclude that this is evidence for hidden heritability due
20 to heterogeneity across the sample population or cohort.

21 Our approach allows us to decompose average heritability levels across historical
22 cohorts and countries into a genetic component that is either ‘universal’ across all environments
23 versus ‘environmentally specific’, enabling a test of whether the same genes are explaining

1 variance in the phenotype to the same extent in different geographical (country) and historical
2 (birth cohort) environments. To test for alternative explanations for heterogeneity across
3 sampling populations, such as genotyping error, we conduct a series of simulation studies to
4 evaluate the role of gene-environment interaction in contrast to alternative explanations. A
5 recent study used bivariate GREML models to investigate genetic heterogeneity in height and
6 BMI between two populations in the US and Europe, providing evidence for homogeneity in
7 both phenotypes.²¹ We expect negligible gene-environment interaction for these
8 anthropometric traits and compare findings for these homogeneous phenotypes to those from
9 our behavioural phenotypes (education, human reproductive behavior) using the same
10 modeling framework.

11 **Material & Methods**

12 **Data**

13 We pooled a series of large datasets consisting of unrelated genotyped men and women
14 (individuals with a >0.05 relatedness as estimated using common SNP markers were removed)
15 from six countries and seven sampling populations in the US (HRS (N=8,146), ARIC
16 (N=6,633)), the Netherlands (LifeLines (N=6,021)), Sweden (STR/SALT (N=6,040)),
17 Australia (QIMR (N=1,167)), Estonia (EGCUT (N=3,722)); and the UK (TwinsUK
18 N=3,333)), for total sample size of N=35,062 (see Text S1 for details on data sources).

19 We used genotype data from all cohorts, imputed to the 1000 genome panel. We then
20 selected HapMap3 SNPs with an imputation score larger than 0.6, excluded SNPs with a
21 missing rate greater than 5%, a lower minor allele frequency than 1% and those which failed
22 the Hardy-Weinberg equilibrium test for a threshold of 10^{-6} . We subsequently applied these

1 criteria again after merging each dataset. We utilized 847,278 SNPs in analyses. The software
2 PLINK²² was used for quality control and merging.

3

4 **Phenotypes**

5 The phenotypes under study are education, human reproductive behavior (number of children
6 ever born (NEB) and age at first birth (AFB)), height, and BMI. We received measures of
7 height and BMI from all cohorts in centimeters and kg/m² respectively or already Z-
8 transformed by sex. For education and human reproductive behavior, we received the
9 phenotypes, which cohorts have used in the respective large-scale GWAS meta-analyses, or
10 constructed them based on raw data and Z-transformed the phenotypes for sex and birth cohorts
11 by dataset.^{4,23}

12 The number of years of education was constructed based on educational categories with
13 the typical years of education in the countries following the standard ISCED scale.^{3,4} The
14 number of children ever born (NEB) measures number of children a woman has given birth to
15 or a man has fathered.² This measure was available in all cohorts, although in ARIC and
16 TwinsUK, only available for women. Information on age at first birth (AFB) was available for
17 all cohorts except for ARIC and HRS. We focus only on individuals who reached the end of
18 their reproductive period of 45 for women and 50 for men (for more details see Text S2).
19 Reproductive phenotypes are frequently recorded, virtually immune to measurement error and
20 used as key parameters for demographic forecasting.²⁴

21

22 **GREML Models**

23 We first describe the baseline GREML model, which assumes the absence of gene-
24 environment interactions. We then extend this model to the GCI-GREML model^{9,20} including

1 genetic relatedness matrices where we stratify data by environments, setting pairwise
 2 relatedness for individuals in different environments to zero.⁹ Doing so allows us to test
 3 whether the pairwise genetic relatedness is a better predictor of pairwise phenotypic similarity
 4 if both individuals live in the same environment, and thus test for gene-environment
 5 interaction. Note that we define the various genetic variance components of the models
 6 explicitly, and will refer to h_{SNP}^2 as the sum of all genetic effects relative to the phenotypic
 7 variance within the respective model specification.

8

9 **Baseline model (GREML)**

10 The genetic component underlying a trait is commonly quantified in terms of SNP-based
 11 heritability as the proportion of the additive genetic variance explained by common SNPs
 12 across the genome over the overall phenotypic variance (σ_Y^2) of the trait: ⁸

$$13 \quad h_{SNP}^2 = \frac{\sigma_G^2}{\sigma_Y^2}$$

14 The phenotypic variance is the sum of additive genetic and environmental variance,
 15 i.e., $\sigma_Y^2 = \sigma_G^2 + \sigma_E^2$, where σ_G^2 is the additive genetic variance explained by all common SNPs
 16 across the genome and σ_E^2 is residual variance. The methods we applied have been detailed
 17 elsewhere.^{8,9,25-27} Briefly, we applied a linear mixed model:

$$18 \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{g} + \mathbf{e}$$

19 where \mathbf{y} is an $N \times 1$ vector of dependent variables, N is the sample size, $\boldsymbol{\beta}$ is a vector for fixed
 20 effects of the M covariates in $N \times M$ matrix \mathbf{X} (including the intercept and potential confounders
 21 such as birth year), \mathbf{g} is the $N \times 1$ vector with each of its elements being the total genetic effect
 22 of all common SNPs for an individual, and \mathbf{e} is an $N \times 1$ vector of residuals. We have $\mathbf{g} \sim N(0, \mathbf{A}\sigma_G^2)$
 23 and $\mathbf{e} \sim N(0, \mathbf{I}\sigma_E^2)$. Hence, the variance matrix \mathbf{V} of the observed phenotypes is:

$$\mathbf{V} = \mathbf{A}\sigma_G^2 + \mathbf{I}\sigma_E^2,$$

To estimate the GRM, 847,278 HapMap3 SNPs were used to capture common genetic variation in the human genome.²⁸ For each individuals (j and k), the corresponding element of the GRM is defined as:

$$A_{jk} = \frac{1}{K} \sum_{i=1}^K \frac{(x_{ij} - 2p_i)(x_{ik} - 2p_i)}{2p_i(1 - p_i)},$$

where x_{ij} denotes the number of copies of the reference allele for the i^{th} SNP for the j^{th} individual and p_i the frequency of the reference allele and K the number of SNPs. If two individuals had a higher genetic relatedness than 0.05, one was excluded from the analyses to avoid bias due to confounding by shared environment amongst close relatives. GCTA was used for the construction of the GRM and GREML analyses.⁹

In the baseline model we apply this approach to the pooled data sources without environmental strata. Hence, the baseline model creates a reference point for SNP-based heritability in the mega-analysis.

Gene \times sampling population ($\mathbf{G}\times\mathbf{P}$) GCI-GREML model

In the case where genetic effects are heterogeneous across sampling populations, SNP-based heritability estimates obtained from the baseline model will be deflated when sampling populations are pooled. We therefore apply a gene \times sampling population model ($\mathbf{G}\times\mathbf{P}$) to simultaneously estimate within- and between variance explanations of common SNPs (see also^{9,20} for GCI-GREML models).

The $\mathbf{G}\times\mathbf{P}$ model jointly model estimates global genetic effects for the outcome variables, effectively between and within samples (σ_G^2) and the averaged additional genetic effects within sampling populations ($\sigma_{\mathbf{G}\times\mathbf{P}}^2$).

1

2

$$\mathbf{V} = \mathbf{A}\sigma_{\mathbf{G}}^2 + \mathbf{A}_{\mathbf{G}\times\mathbf{P}}\sigma_{\mathbf{G}\times\mathbf{P}}^2 + \mathbf{I}\sigma_{\mathbf{E}}^2$$

3

where \mathbf{A} is the genetic relatedness matrix and $\mathbf{A}_{\mathbf{G}\times\mathbf{P}}$ is a matrix only with values for pairs of

4

individuals within Populations 1–7:

5

6

$$\mathbf{A} = \begin{bmatrix} A_{p1p1} & A_{p2p1} & A_{p3p1} & A_{p4p1} & A_{p5p1} & A_{p6p1} & A_{p7p1} \\ A_{p1p2} & A_{p2p2} & A_{p3p2} & A_{p4p2} & A_{p5p2} & A_{p6p2} & A_{p7p2} \\ A_{p1p3} & A_{p2p3} & A_{p3p3} & A_{p4p3} & A_{p5p3} & A_{p6p3} & A_{p7p3} \\ A_{p1p4} & A_{p2p4} & A_{p3p4} & A_{p4p4} & A_{p5p4} & A_{p6p4} & A_{p7p4} \\ A_{p1p5} & A_{p2p5} & A_{p3p5} & A_{p4p5} & A_{p5p5} & A_{p6p5} & A_{p7p5} \\ A_{p1p6} & A_{p2p6} & A_{p3p6} & A_{p4p6} & A_{p5p6} & A_{p6p6} & A_{p7p6} \\ A_{p1p7} & A_{p2p7} & A_{p3p7} & A_{p4p7} & A_{p5p7} & A_{p6p7} & A_{p7p7} \end{bmatrix}$$

7

8

$$\mathbf{A}_{\mathbf{G}\times\mathbf{P}} = \begin{bmatrix} A_{p1p1} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & A_{p2p2} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & A_{p3p3} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & A_{p4p4} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & A_{p5p5} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & A_{p6p6} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & A_{p7p7} \end{bmatrix}$$

9

10

11

The sum of both variance components ($\sigma_{\mathbf{G}}^2 + \sigma_{\mathbf{G}\times\mathbf{P}}^2$) are therefore expected to

12

correspond with the results of a meta-analysis of the sample-specific h_{SNP}^2 of sufficient sample

13

size. We quantify the hidden heritability due to heterogeneity as the discrepancy between

14

$h_{\text{SNP}}^2 = \frac{\sigma_{\mathbf{G}}^2}{\sigma_{\mathbf{Y}}^2}$ from the baseline model and $h_{\text{SNP}}^2 = \frac{\sigma_{\mathbf{G}}^2 + \sigma_{\mathbf{G}\times\mathbf{P}}^2}{\sigma_{\mathbf{Y}}^2}$ from the GxP model.

15

16

Gene \times demographic birth cohort (G \times C) GCI-GREML model

1 We are likewise interested in gene-environment interaction across birth cohorts. Fertility
 2 behavior and educational attainment have dramatically changed during the 20th century.^{18,29}
 3 Figure 1 shows the trends in age at first birth (AFB) during the 20th century for the countries
 4 in our study (see Text S3 for details on the data sources). We see the well-established U-shaped
 5 pattern of a falling AFB in the first half of the 20th century followed by an upturn in the trend
 6 of AFB towards older ages. This widespread fertility postponement¹⁸ – referred to as the
 7 Second Demographic Transition³⁰ – was related to the spread of effective contraception, a
 8 drop in the NEB, changes in the economic need for children and female educational
 9 expansion.^{18,31}

10 Environmental changes occurred at different periods in each country, with Australia
 11 having the earliest onset of fertility postponement (1939) and Estonia having the latest due to
 12 post-socialist transitions (1962; see Table S1 for all turning points and details). To test for gene-
 13 environment interaction, we grouped the birth cohorts into environmentally homogeneous
 14 conditions by those born before and after each country-specific fertility postponement turning
 15 point. To investigate the moderating effect of turning points, we follow the previous modeling
 16 strategy, but divide individuals into these turning point birth cohorts.

17 (Figure 1 here)

18 The $\mathbf{G} \times \mathbf{C}$ model is a joint model estimating the universal genetic effects for the traits,
 19 effectively between and within samples ($\sigma_{\mathbf{G}}^2$) and the averaged additional genetic effects
 20 within defined birth cohorts ($\sigma_{\mathbf{G} \times \mathbf{C}}^2$):

$$21 \quad \mathbf{V} = \mathbf{A}\sigma_{\mathbf{G}}^2 + \mathbf{A}_{\mathbf{G} \times \mathbf{C}}\sigma_{\mathbf{G} \times \mathbf{C}}^2 + \mathbf{I}\sigma_{\mathbf{E}}^2$$

22 where \mathbf{A} is the genetic relatedness matrix and $\mathbf{A}_{\mathbf{G} \times \mathbf{C}}$ is a matrix only with values for pairs of
 23 individuals within the same demographic birth Cohorts $c_1 - c_2$:

$$\mathbf{A}_{\mathbf{G} \times \mathbf{C}} = \begin{bmatrix} \mathbf{A}_{c_1 c_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{c_2 c_2} \end{bmatrix}$$

Genes × Population × Demographic birth cohorts (G×P×C) GCI-GREML model

In the **G×P×C** model, we included both interaction terms mentioned above and an additional interaction term $\mathbf{A}_{\mathbf{G} \times \mathbf{P} \times \mathbf{C}}$ which is equal to zero for all pairs of individuals living in different time periods or in different cohorts represented by:

$$\mathbf{V} = \mathbf{A}\sigma_{\mathbf{G}}^2 + \mathbf{A}_{\mathbf{G} \times \mathbf{P}}\sigma_{\mathbf{G} \times \mathbf{P}}^2 + \mathbf{A}_{\mathbf{G} \times \mathbf{C}}\sigma_{\mathbf{G} \times \mathbf{C}}^2 + \mathbf{A}_{\mathbf{G} \times \mathbf{P} \times \mathbf{C}}\sigma_{\mathbf{G} \times \mathbf{P} \times \mathbf{C}}^2 + \mathbf{I}\sigma_{\mathbf{E}}^2$$

where \mathbf{A} is the genetic relatedness matrix, $\mathbf{A}_{\mathbf{G} \times \mathbf{P}}$ is a matrix only with non-zero values for pairs of individuals within populations from the **G×P** Model, $\mathbf{A}_{\mathbf{G} \times \mathbf{C}}$ is a matrix only with non-zero values for pairs of individuals within the same demographic periods from the **G×C** Model, and $\mathbf{A}_{\mathbf{G} \times \mathbf{P} \times \mathbf{C}}$ is a matrix only with values for pairs of individuals with both the same demographic periods and the same populations.

Control variables

All phenotypes have been Z-transformed by sampling population, birth year and sex. We furthermore added fixed effects for sex, birth year, sampling population (with reference category Lifelines, the Dutch dataset) and the first 20 principal components calculated from the GRM across all populations to account for population stratification.³² For the interaction model with birth cohorts, we included an additional fixed effect for the respective birth cohort turning point. In the **G×P×C** model, we additionally controlled for the interactions between the respective sampling population and the birth cohort division.

Model-fitting approach

1 The variance components are estimated using GREML estimation. When comparing the
2 respective model specifications, to determine the best-fitting model, we rely on a model-fitting
3 approach that compares the full model with reduced models that constrain specific effects to
4 be zero. Since the models are nested, we perform likelihood-ratio tests and prefer the more
5 parsimonious models if there is no significant loss in model fit (where the test statistic is
6 distributed as a mixture of chi-squared with a probability of 0.5 and 0⁹; p-values from these
7 tests are provided in Tables S3-7).⁹ This strategy is also robust against the violation of the
8 assumption of requiring a normal distribution of the dependent variable – as for example in the
9 case of NEB (number of children ever born).³³

10

11 **Simulation Study**

12 We conducted a series of simulation studies to illustrate how our models interpret gene-
13 environment interaction and to evaluate the role of potential alternative sources of
14 heterogeneity in our data. All simulation studies are detailed in Text S4 (for the theory behind
15 them see ²⁰). First, we were interested in how the model construes heterogeneity in heritability
16 levels across populations. Since heritability is a ratio of the proportion of total phenotypic
17 variance that is attributable to additive genetic effects, differences in the residual variance for
18 example due to heterogeneous phenotypic measurement error can lead to different levels of
19 heritability across populations, even though genetic effects are perfectly correlated. In contrast
20 to twin studies, we are not interested in comparing levels of heritability across populations, but
21 in the question of whether genes have the same effect on the phenotype across environments.
22 We thus decompose the heritability in the pooled data into additive genetic variance, both
23 within and between environments.

1 In simple terms, we simulated phenotypes without gene-environment interaction across
2 sampling populations and with gene-environment interaction across sampling populations
3 based on 5000 SNPs that were in approximate linkage equilibrium (pairwise r^2 between SNPs
4 below 0.05) and repeated this across 50 replications. First, to test for a model without gene-
5 environment interaction, we set h_{SNP}^2 of the trait to 0.50 and the genetic correlations across
6 environments to 1 (Text S4 Sim 1). Second, we repeated the simulations with varying residual
7 phenotypic variance across populations³⁴, resulting in simulated h_{SNP}^2 between 0.25–0.625,
8 but still with a genetic correlation of 1 across populations (Text S4 Sim 2). Third, to illustrate
9 weak levels of gene-environment interaction, we simulated h_{SNP}^2 to be 0.50 and the genetic
10 correlations of traits across populations to be 0.80 (Text S4 Sim 3). Finally, to illustrate stronger
11 gene-environment interaction, we simulated h_{SNP}^2 to 0.50 and the genetic correlations of traits
12 across populations to 0.50 (Text S4 Sim 4).

13 The stacked bars in Figure 2 depict the average estimates of the four types of
14 simulations for the simulated 50 phenotypes for the baseline model and the $\mathbf{G} \times \mathbf{P}$ model
15 (individual estimates are presented as black dots for the full model and stripes in the bars
16 represent variance components). Examining the first model (Sim 1) assumed no gene-
17 environment interaction by sampling populations and thus homogeneous heritability, h_{SNP}^2 as
18 $\sigma_{\mathbf{G}}^2/\sigma_{\mathbf{Y}}^2$ (blue bar) is estimated at 0.324 and therefore around three fifths of the simulated
19 heritability of 0.50 since the GRM is based not only on quantitative trait loci. Central to our
20 approach is that for the phenotypes with no $\mathbf{G} \times \mathbf{P}$ interaction, the variance explanation that is
21 effective both within and between populations ($\sigma_{\mathbf{G}}^2/\sigma_{\mathbf{Y}}^2$) is nearly identical to the baseline model
22 (0.318). The gene-environment interaction term ($\sigma_{\mathbf{G} \times \mathbf{P}}^2/\sigma_{\mathbf{Y}}^2$) estimates a small additional
23 explanation of variance within populations of on average 0.026, with the full model estimate

1 of h_{SNP}^2 within populations at 0.344 $\left(= \frac{\sigma_{\text{G}}^2 + \sigma_{\text{G}\times\text{P}}^2}{\sigma_{\text{Y}}^2}\right)$. Importantly, the same holds if we simulate
2 differences in h_{SNP}^2 across populations due to varying residual variance. Sim 2 in Figure 2
3 shows an average h_{SNP}^2 of 0.205 and the $\mathbf{G}\times\mathbf{P}$ interaction model estimates of ‘universal’ genetic
4 variance $(\sigma_{\text{G}}^2/\sigma_{\text{Y}}^2)$ of 0.200, with a gene-environment interaction term $(\sigma_{\text{G}\times\text{P}}^2/\sigma_{\text{Y}}^2)$ of 0.0217. We
5 therefore conclude that the model does not interpret heterogeneity in heritability levels due to
6 differences in the residual variance as gene-environment interaction.

7 (Figure 2 here)

8 Sim 3 and 4 in Figure 2 depict how gene-environment interaction across sampling
9 populations affects model estimates in scenarios of cross population genetic correlations of
10 0.80 (weak) and 0.50 (strong) gene-environment interaction respectively, but the same
11 population specific h_{SNP}^2 of 0.050 as in Sim 1. First, we observe that h_{SNPs}^2 in the baseline
12 models are deflated in the pooled data $\left(\frac{\sigma_{\text{G}}^2}{\sigma_{\text{Y}}^2} = 0.261 \text{ and } 0.105\right)$ and therefore only capture
13 around four-fifths and one-third of the estimates in the absence of $\mathbf{G}\times\mathbf{P}$. Second, when taking
14 $\mathbf{G}\times\mathbf{P}$ into account, the full model estimate reaches the same level as the baseline model in the
15 absence of $\mathbf{G}\times\mathbf{P}$ $\left(\frac{\sigma_{\text{G}}^2 + \sigma_{\text{G}\times\text{P}}^2}{\sigma_{\text{Y}}^2} = 0.328 \text{ and } 0.315\right)$ due to a larger fraction of genetic variance
16 explained within populations $\left(\frac{\sigma_{\text{G}\times\text{P}}^2}{\sigma_{\text{Y}}^2} = 0.082 \text{ and } 0.256\right)$ and do not appear to be inflated
17 whatsoever. Third, the genetic variance explained effectively within and between populations
18 in the $\mathbf{G}\times\mathbf{P}$ model is even smaller than in the baseline model $\left(\frac{\sigma_{\text{G}}^2}{\sigma_{\text{Y}}^2} = 0.246 \text{ and } 0.059\right)$.
19 Therefore, while in the case of a genetic correlation of 0.5 across populations, within
20 population estimates of h_{SNP}^2 capture around one third of the overall heritability; the shared
21 genetic variance explanation across populations would be only around 19% $(=0.059/0.315)$ of
22 this value.

1 Based on the findings from Sim 4 for example, we would expect that in the case of
2 meta-analyses of population specific GWAS on the gene-environment interaction phenotypes,
3 that genome-wide significant SNPs could explain only up to 10% of the variance while h_{SNP}^2
4 of within populations could explain on average 32%. Around 68% of h_{SNP}^2 $((1-10/32)*100)$
5 would therefore be ‘hidden’ in the mega-analysis due to heterogeneity and in this case due to
6 gene-environment interaction.

7 Figure 3 shows hidden heritability estimates for the simulations without gene-
8 environment interaction (Sim 1) and with gene-environment interaction (Sim 3 and Sim 4). We
9 were furthermore interested to what extent genetic heterogeneity across populations such as
10 differences in genetic measurement, in linkage disequilibrium across sampling populations, or
11 heterogeneous imputation quality across population can lead to observed heterogeneity or
12 deflate h_{SNP}^2 in pooled data sources. To investigate this we removed the 5,000 causal SNPs
13 from the genetic data, which was the basis of how we simulated the phenotypes. We then re-
14 estimated the GRM and repeated the analyses on Sim 1, of phenotypes without gene-
15 environment interaction and homogeneous heritability across populations (depicted in Figure
16 3 as Sim 1 LD). If the causal SNPs are removed, estimates are based on correlated SNPs, which
17 are in linkage disequilibrium (LD). To the extent that the structure in the genetic data we use
18 is heterogeneous across populations due to the aforementioned reasons, we can expect that our
19 models interpret it as heterogeneous genetic effects resulting in hidden heritability.

20 In Figure 3, we see that hidden heritability is estimated to be around 68% for a genetic
21 correlation of 0.50, around 20% for a genetic correlation of 0.80 and around 5% for the model
22 without gene-environment interaction as well as a model based on SNPs in LD with the causal
23 SNPs. This allows us to draw two conclusions. First, in the complete absence of gene-

1 environment interaction (Sim 1), our models interpret on average across 50 simulations that
2 5% of the heritability in the $\mathbf{G}\times\mathbf{P}$ model is hidden in a standard model with a statistically
3 significant $\mathbf{G}\times\mathbf{P}$ term in 10 simulation studies ($10/50 = 20\%$; not listed) at the 5%-level. This
4 is important to keep in mind when analyzing our phenotypes of interest. To evaluate phenotype
5 specific model inflations, we conducted complementary permutation analyses generating a
6 matrix with randomly stratified environments to see how estimates are inflated in the real data
7 for specific phenotypes. This will be reported when discussing the findings. Second, we find
8 no difference in inflation between the simulations including and excluding causal SNPs (Sim
9 1 LD and Sim 1). We conclude from this that heterogeneity in the genetic structure of the
10 populations does not affect our interpretation of gene-environment interaction in comparison
11 to the standard model. This is likely due to the fact that we only look at common SNPs and
12 applied rigorous quality control. To investigate whether gene-environment interaction is
13 present for education and human reproductive behavior, we estimate the above models as well
14 as $\mathbf{G}\times\mathbf{C}$ and $\mathbf{G}\times\mathbf{P}\times\mathbf{C}$ models to these phenotypes in seven sampling populations.

15

16 **Sex differences**

17 Previous studies find no evidence for gene-sex interaction of common genetic effects on BMI,
18 height³⁵ and also human reproductive behavior². We also tested for G x Sex interaction within
19 sampling populations, as:

$$20 \quad \mathbf{V} = \mathbf{A}_{\mathbf{G}\times\mathbf{P}}\sigma_{\mathbf{G}\times\mathbf{P}}^2 + \mathbf{A}_{\mathbf{G}\times\mathbf{P}\times\text{sex}}\sigma_{\mathbf{G}\times\mathbf{P}\times\text{sex}}^2 + \mathbf{I}\sigma_{\mathbf{E}}^2$$

21 where $\mathbf{A}_{\mathbf{G}\times\mathbf{P}}$ is the genetic relatedness matrix only with values for pairs of individuals within
22 the same population and $\mathbf{A}_{\mathbf{G}\times\mathbf{P}\times\text{sex}}$ is a matrix with only values for pairs of individuals of the
23 same sex and same sampling population.

1 (Figure 3 here)

2 Decomposing the genetic variance of all five phenotypes, height, BMI, education,
3 number of children ever born (NEB) and age at first birth (AFB) into within population effects
4 shared between sexes ($\sigma_{\mathbf{G}\times\mathbf{P}}^2$) and the averaged additional genetic effects within sexes
5 ($\sigma_{\mathbf{G}\times\mathbf{P}\times\mathbf{sex}}^2$), we find no evidence for sex-specific effects ($\sigma_{\mathbf{G}\times\mathbf{P}\times\mathbf{sex}}^2$) for education (p-value
6 0.49), AFB (p-value 0.5), NEB (p-value 0.41) or height (p-value 0.5). Only for BMI do we find
7 evidence of around a 3% sex-specific variance explanation (p-value 0.046; for full results see
8 Table S2). Given that we focus on education and reproductive behavior, we applied all models
9 to pooled data including both sexes, keeping in mind the findings for BMI.

10

11 **Results**

12 **SNP-based heritability across model specifications by phenotypes**

13 When we ignore environmental differences, h^2_{SNP} in the standard GREML model (G) is
14 significant for all phenotypes, but at different levels (Figure 4 and Tables S3-S7 for full model
15 estimates). For height, h^2_{SNP} is estimated as 0.40 (SE 0.01), meaning that 40% of the variance
16 in height can be attributed to common additive genetic effects. h^2_{SNP} is smaller for BMI (0.17
17 SE 0.01) and years of education (0.16 SE 0.01) and low for both reproductive behavior
18 outcomes, NEB (0.03 SE 0.01) and AFB (0.08 0.02) respectively.

19 More importantly, however, for our question, h^2_{SNP} in all phenotypes increases if we
20 include stratified GRMs in addition to the baseline GRM (e.g., yielding the $\mathbf{G}\times\mathbf{C}$ model when
21 stratifying by birth cohorts, the $\mathbf{G}\times\mathbf{P}$ model when stratifying by sampling populations, and the
22 $\mathbf{G}\times\mathbf{P}\times\mathbf{C}$ model when stratifying by both). Particularly for the complex behavioral outcomes of
23 education and reproductive behavior, the increase is substantial. For education, h^2_{SNP} increases

1 by 80% (up to 0.28 SE 0.03) in the $\mathbf{G}\times\mathbf{P}\times\mathbf{C}$ model compared to the standard GREML model
2 (\mathbf{G}). For AFB, the increase is 60% (0.13 SE 0.04) and for NEB it is as high as 342% (0.13 SE
3 0.03). In contrast, the increase in the full $\mathbf{G}\times\mathbf{P}\times\mathbf{C}$ model was considerably smaller at 12% (0.44
4 SE 0.03) for height and 30% (0.22 SE 0.03) for BMI.

5

6 **Best model by phenotype**

7 Based on likelihood ratio tests, we identified the best fitting and parsimonious model (in Figure
8 4 marked as BM; for full results see Table S8). For height, the best fitting model includes no
9 gene-environment interaction and therefore corroborates previous findings from the literature.

10 ³⁵

11 (Figure 4 here)

12 For BMI, and the reproductive phenotypes of AFB and NEB, the $\mathbf{G}\times\mathbf{P}$ specification shows the
13 best model fit. This indicates significant gene-environment interaction across sampling
14 populations, while there is no evidence for an interaction by birth cohort. For BMI, additive
15 SNP variance, or effectively between and within populations (i.e., the blue column that
16 assumes it is effective across the defined environments or ‘universal’ respectively; $\sigma_{\mathbf{G}}^2/\sigma_{\mathbf{Y}}^2$),
17 16% of the variance in the phenotype and an additional 5% can be explained on average within
18 populations ($\sigma_{\mathbf{G}\times\mathbf{P}}^2/\sigma_{\mathbf{Y}}^2$, green column). For AFB, around 6% of the variance can be explained
19 by universal genetic effects while 7% are environmentally specific, and for NEB only 1% of
20 the variance can be explained between populations, with 12% within them. Finally, for
21 education, the best-fitting model ($\mathbf{G}\times\mathbf{P}\times\mathbf{C}$) implies that both sampling population and birth
22 cohort moderate genetic effects from the whole genome and that there are genetic effects

1 unique to sampling populations within the defined birth cohorts. In contrast to reproductive
2 behavior, however, 12% of the overall variance can still be explained by additive common
3 genetic effects even between populations. Additionally, there is 2% variance explained within
4 birth cohorts ($\sigma_{\mathbf{G} \times \mathbf{C}}^2 / \sigma_{\mathbf{Y}}^2$, red column), 6% within populations and 8% which is unique within
5 populations and birth cohorts ($\sigma_{\mathbf{G} \times \mathbf{P} \times \mathbf{C}}^2 / \sigma_{\mathbf{Y}}^2$, orange column).

6 **Quantifying ‘universal effects’ and ‘hidden heritability’ due to heterogeneity**

7 Figure 5 visualizes: (i) the ‘universal effects’ or ratio for genetic variance captured by the
8 normal GRM in the best fitting model (i.e., blue column, $\sigma_{\mathbf{G}}^2 / \sigma_{\mathbf{Y}}^2$ in the model with the best fit)
9 and the total h_{SNP}^2 (i.e., across all genetic components in the best fitting model). It also shows
10 (ii) in red the ‘hidden heritability’ due to gene-environment interaction (i.e., the differences in
11 total h_{SNP}^2 between the best fitting model and the baseline model, divided by the total h_{SNP}^2 of
12 the best fitting model) for all phenotypes.

13 (Figure 5 here)

14
15 The Figure illustrates hidden heritability due to gene-environment interaction
16 particularly for the complex phenotypes we are most interested in, namely: education, and the
17 reproductive outcomes of AFB and NEB. For education, only 55% of h_{SNP}^2 in the best fitting
18 model is ‘universal’ or effectively both within and between environments. A standard GREML
19 model (G) would only capture around 63% of h_{SNP}^2 in the best fitting model resulting in 37%
20 hidden heritability. For reproductive behavior, this becomes even stronger. For NEB only 6%
21 of h_{SNP}^2 of h_{SNP}^2 of the best fitting model is universal, with 75% hidden in the baseline model.
22 For AFB, 45% of h_{SNP}^2 is universal with around 40% of the h_{SNP}^2 hidden in the baseline model.
23 For height, in contrast, we see that the h_{SNP}^2 in the best fitting model is effectively between
24 environments and we find no evidence for hidden heritability. For BMI, around 75% of h_{SNP}^2

1 in the best fitting model is effectively between and within environments (i.e., universal). The
2 standard GREML model (G) for BMI thus captures 80% of h_{SNP}^2 from the best fitting model
3 with 20% hidden heritability.

4

5 **Discussion**

6 Using whole-genome data from seven populations, we demonstrate heterogeneity in genetic
7 effects across national populations and birth cohorts for educational attainment and human
8 reproductive behavior in a mega-analysis framework. Our findings imply substantial ‘hidden
9 heritability’ due to heterogeneity for educational attainment (37%) and reproductive behavior
10 (40% for AFB and 75% for NEB) in the cohorts under study. Comparative analysis with
11 anthropometric traits (height and BMI) corroborates previous findings from whole-genome
12 methods of a more homogeneous genetic architecture of these phenotypes across environments
13 (while for BMI GWA studies also find evidence for gene-environment interaction across birth
14 cohorts in the HRS ^{36,37}).

15 Our findings indicate that the lower predictive power of polygenic scores from large
16 GWA studies compared to SNP-based heritability on single or very few populations partly
17 reflects the fact that genetic effects are (to some extent) not universal but rather specific to data
18 sources for these complex traits. Estimates are well in line with the 36-38% loss in polygenic
19 score R^2 across data sets reported for education.³⁸ They demonstrate therefore that the reference
20 SNP-based heritability for the predictive power of polygenic scores obtained from the GWAS
21 meta-analyses amongst several populations is smaller than SNP-based heritability obtained
22 from single populations. While the need for statistical power often still necessitates large-scale
23 GWAS meta-analysis combining multiple and diverse data sources, our findings also suggests

1 that large homogeneous data sources such as the UK Biobank with around 500,000 genotyped
2 individuals may trigger genetic discovery for behavioral outcomes. Drawing conclusions or
3 making predictions out of one discovery sample alone, however, may be inaccurate, since SNPs
4 may have different effects in different samples, or the phenotype may reflect different
5 behavioral aspects.

6 Complementary simulation studies corroborate the interpretation that our findings are
7 mainly driven by gene-environment interaction in contrast to heterogeneity in residual
8 environmental variance – including measurement error – or genetic heterogeneity (e.g.,
9 genotyping platform, genetic architecture, imputation quality) across the data sources we
10 pooled. When applying our models to simulated phenotypes without gene-environment
11 interaction but rather to different levels of heritability due to varying residual variance, we find
12 no systematic inflation of the **GxP** component in our models. Furthermore, both estimates
13 including and excluding the causal 5000 SNPs result in nearly identical estimates. In the total
14 absence of gene-environment interaction, estimates show a slight inflation in the **GxP** model
15 (5%). First, the substantial findings of hidden heritability between 40–75% for behavioral
16 phenotypes largely exceeds this potential inflation, corresponding with simulations of a genetic
17 correlation between 0.5–0.8 across populations for the behavioral phenotypes. Second, we
18 conducted permutation analyses, generating a random gene-environment interaction, not
19 stratifying by population or birth cohorts. Here we found no inflation for age at first birth by a
20 randomly generated matrix included in the models ($\sigma_{\mathbf{G}\times\mathbf{P}}^2$ 0.000001, SE 0.03, p-value 0.50), nor
21 for number of children ever born ($\sigma_{\mathbf{G}\times\mathbf{P}}^2$ 0.003, SE 0.02, p-value 0.43) nor education ($\sigma_{\mathbf{G}\times\mathbf{P}}^2$
22 0.000001, SE 0.02, p-value 0.50; not listed). It remains vital to conclude that although the
23 estimates of hidden heritability provided in our study are the first estimates of its kind in a

1 single design – in contrast to comparing GWAS and whole-genome methods – estimates do
2 not represent generalizable values of hidden heritability for these traits. The estimates are for
3 illustrative purposes as they are also dependent on the number of cohorts combined for a study,
4 the respective level of heterogeneity across them and they might be slightly inflated.

5 Contrary to our expectations, we did not find any evidence for gene-environment
6 interaction across birth cohorts for human reproductive behavior. This is particularly surprising
7 since across time there have been substantial environmental changes such as the introduction
8 of effective contraception, social norms around the timing of childbearing and educational
9 expansion – all factors which strongly modifies reproductive behavior.¹⁸ In contrast, we find
10 cohort specific genetic effects on educational attainment. This contributes to solving the puzzle
11 of missing heritability in educational attainment, since twin studies with higher heritability
12 estimates are also conducted within homogeneous birth cohorts.

13 Our findings expose the challenges in detecting genetic variants associated with human
14 reproductive behavior or other complex phenotypes in GWAS meta-analyses of multiple
15 cohorts. First, SNP-based heritability within populations is comparably small and second, we
16 find limited evidence that genetic effects underlying reproductive behavior in one country
17 predicts the underlying behavior in another. Our findings likely reflect the interrelated
18 behavioral nature of reproduction and education, which appears to be more sensitive to cultural
19 and societal heterogeneity than for example anthropometric traits such as height or BMI. It has
20 also been shown that pleiotropic genes affecting age at first birth and schizophrenia have
21 different effects across populations.³⁹ Recently, social scientists have made considerable
22 efforts to integrate molecular genetics into their research.²⁻⁴ When considering the highly
23 socially- and biologically-related phenotype of reproductive behavior outcomes,
24 environmental factors are critical in understanding how genetic factors are modified in relation

1 to fecundity and infertility.

2 Finally, our study also has several important limitations. First, it is possible that
3 heterogeneity in the phenotypic measures influences the patterns we observed. While we find
4 no evidence that our models interpret changing relative environmental contributions to trait
5 variation as gene-environment interaction, we cannot rule out the possibility that the trait
6 definitions differ across environments. We consider this as minor issue for reproductive
7 behavior. While measures are not perfectly harmonized across birth cohorts (for e.g., some
8 questionnaires for example explicitly ask for number of still-births and others do not), in
9 LifeLines and TwinsUK, we compared the live birth measures with number of children ever
10 born and, as expected, given the low mortality rate in both populations, less than 0.2% of the
11 children had not reached reproductive age. Moreover, the correlation of number of children
12 ever born and number of children reaching reproductive age was 0.98. We therefore do not
13 expect a large bias due to the exclusion of stillbirths in some countries (for details see Text S2).
14 Nevertheless, we cannot reject the possibility that heterogeneity in the measure of education
15 remains even after homogenizing it with the standard ISCED scale. In this case, we would
16 argue that large parts of the gene-environment interaction pattern we observe for education are
17 due to interaction within populations by birth cohorts where we hypothetically have
18 homogeneous measures. Furthermore, different cross-national definitions of education
19 represent a case of gene-environment interaction. Finally, our statistical findings of
20 heterogeneity are of major importance in shaping our expectations about the ability to locate
21 genetic loci associated with education in GWAS meta-analyses despite their causal
22 mechanisms.

23 Second, notwithstanding the fact that our simulation studies show no inflation of hidden
24 heritability due to differences in the genetic structure across populations, it is plausible that

1 empirical phenotypes are heterogeneous in reference to rare genetic variants, which are not
2 considered in our models and not present in our data. This is an issue demanding further
3 consideration in future research. We are suitably cautious that part of hidden heritability in our
4 models might be driven by rare, population-specific variants. Previous studies of height and
5 BMI show that rare variants explain a significant part of phenotypic variance,¹⁵ while our
6 models show the least heterogeneity across populations for these phenotypes.

7 Third, the models we apply average within environmental effects across populations.
8 An optimal study design would be a multivariate genetic modeling approach, which estimates
9 SNP-based heritability for each population and the genetic correlations across them. This
10 approach, however, is feasible for traits with strong or moderate heritability such as height and
11 BMI,²¹ but lack statistical power²⁷ for phenotypes with small SNP-based heritability such as
12 reproductive behavior¹⁰ in the current samples. The models we propose allow us to investigate
13 and compare gene-environment interaction across a range of phenotypes. Multivariate models
14 may become feasible in the future with larger homogeneous data sources, and will also enable
15 us to disentangle shared genetic effects across these phenotypes.^{2,40,41}

16 Finally, in the current modeling approach, we cannot include childless individuals in
17 the modeling of AFB, and future research in quantitative genetics may aim to integrate
18 censored information in their modeling approaches, which is standard in demographic research
19 (for further discussion see ^{10,42,43}).

20 In conclusion, our study provides initial evidence that there may be substantial
21 implications for the estimates of genetic effects due to the extreme heterogeneity of combining
22 multiple samples across vast historical time periods and populations to conduct large-scale
23 GWA studies. We likewise uncover challenges for investigations into the genetic architecture
24 of human reproductive behavior and education and suggest that gene-environment interaction

1 is the main driver of heterogeneity across populations. These challenges can only be overcome
2 by interdisciplinary work between both geneticists and social scientists using ever-larger
3 datasets, with combined information and substantive knowledge of complex phenotypes and
4 environmental conditions.^{44,45}

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

References

1. Visscher, P. M., Brown, M. A., McCarthy, M. I. & Yang, J. Five years of GWAS discovery. *Am. J. Hum. Genet.* **90**, 7–24 (2012).
2. Barban, N. *et al.* Genome-wide analysis identifies 12 loci influencing human reproductive behavior. *Nat. Genet.* 10.1038/ng.3698 (2016). doi:10.1038/ng.3698
3. Rietveld, C. A. *et al.* GWAS of 126,559 Individuals Identifies Genetic Variants Associated with Educational Attainment. *Science (80-.)*. **340**, 1467–1471 (2013).
4. Okbay, A. *et al.* Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* **533**, 539–542 (2016).
5. Polderman, T. J. C. *et al.* Meta-analysis of the heritability of human traits based on fifty years of twin studies. *Nat. Genet.* **47**, 702–709 (2015).
6. Mills, M. C. & Tropf, F. C. The Biodemography of Fertility: A Review and Future Research Frontiers. *Kolner Z. Soz. Sozpsychol.* **55**, 397–424 (2016).
7. Branigan, A. R., McCallum, K. J. & Freese, J. Variation in the heritability of educational attainment: An international meta-analysis. *Soc. forces* **92**, 109–140 (2013).
8. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* **42**, 565–569 (2010).
9. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
10. Tropf, F. C. *et al.* Human fertility, molecular genetics, and natural selection in modern societies. *PLoS One* **10**, e0126821 (2015).
11. Witte, J. S., Visscher, P. M. & Wray, N. R. The contribution of genetic variants to disease depends on the ruler. *Nat. Rev. Genet.* **15**, 765–776 (2014).
12. Manolio, T. A. *et al.* Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009).
13. Zhu, Z. *et al.* Dominance genetic variation contributes little to the missing heritability for human complex traits. *Am. J. Hum. Genet.* **96**, 377–385. (2015).
14. Zuk, O. & Hechter, E. The mystery of missing heritability: Genetic interactions create phantom heritability. *Proc. Natl. Acad. Sci.* **109**, 1193–1198. (2012).
15. Yang, J. *et al.* Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat. Genet.* **47**, 1114–1120 (2015).
16. Felson, J. What can we learn from twin studies? A comprehensive evaluation of the equal environments assumption. *Soc. Sci. Res.* **43**, 184–199 (2014).
17. Wray, N. R. *et al.* Pitfalls of predicting complex traits from SNPs. *Nat. Rev. Genet.* **14**, 507–15 (2013).

- 1 18. Mills, M. C., Rindfuss, R. R., McDonald, P. & te Velde, E. Why do people postpone
2 parenthood? Reasons and social policy incentives. *Hum. Reprod. Update* **17**, 848–860
3 (2011).
- 4 19. de Vlaming, R. *et al.* Meta-GWAS Accuracy and Power (MetaGAP) Calculator Shows
5 that Hiding Heritability Is Partially Due to Imperfect Genetic Correlations across
6 Studies. *PLOS Genet.* **13**, e1006495 (2017).
- 7 20. Robinson, M. R. *et al.* Genotype-covariate interaction effects and the heritability of
8 adult body mass index. *Nat. Genet.*
- 9 21. Yang, J. *et al.* Genome-wide genetic homogeneity between sexes and populations for
10 human height and body mass index. *Hum. Mol. Genet.* **24**, 7445–7449 (2015).
- 11 22. Purcell, S. M. *et al.* PLINK: a tool set for whole-genome association and population-
12 based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 13 23. Barban, N. *et al.* Genome-wide analysis identifies 12 loci influencing human
14 reproductive behavior. *Nat. Genet.* 10.1038/ng.3698 (2016). doi:10.1038/ng.3698
- 15 24. Barban, N, et al. Mills, M. C. Large-scale genomic analysis identifies 12 loci
16 harbouring genes for human reproductive behaviour and infertility. *under Rev.* (2016).
- 17 25. Lee, S. H., Yang, J., Goddard, M. E., Visscher, P. M. & Wray, N. R. Estimation of
18 pleiotropy between complex diseases using single-nucleotide polymorphism-derived
19 genomic relationships and restricted maximum likelihood. *Bioinformatics* **28**, 2540–
20 2542 (2012).
- 21 26. Visscher, P. M., Yang, J. & Goddard, M. E. A commentary on ‘common SNPs explain
22 a large proportion of the heritability for human height’ by Yang et al.(2010). *Twin Res.*
23 *Hum. Genet.* **13**, 517–524 (2010).
- 24 27. Visscher, P. M. *et al.* Statistical power to detect genetic (co) variance of complex traits
25 using SNP data in unrelated samples. *PLoS Genet.* **10**, e1004269 (2014).
- 26 28. Consortium, I. H. 3. Integrating common and rare genetic variation in diverse human
27 populations. *Nature* **467**, 52–58. (2010).
- 28 29. Balbo, N., Billari, F. C. & Mills, M. C. Fertility in advanced societies: A review of
29 research. *Eur. J. Popul. Eur. Démographie* **29**, 1–38 (2013).
- 30 30. de Kaa, D. J. Van. Europe’s second demographic transition. *Popul. Bull.* **42**, 1–59
31 (1987).
- 32 31. Sobotka, T. Is Lowest-Low Fertility in Europe Explained by the Postponement of
33 Childbearing? *Popul. Dev. Rev.* **30**, 195–220 (2004).
- 34 32. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-
35 wide association studies. *Nat. Genet.* **38**, 904–909 (2006).
- 36 33. Snijders, T. A. B. *Multilevel analysis.* (Springer, 2011).
- 37 34. Domingue, B. W. *et al.* Genome-Wide Estimates of Heritability for Social
38 Demographic Outcomes. *Biodemography Soc. Biol.* **62**, 1–18 (2016).
- 39 35. Yang, J. et al. Genome-wide genetic homogeneity between sexes and populations for
40 human height and body mass index. *Hum. Mol. Genet.* **24**, 7445–7449 (2015).
- 41 36. Conley, D., Laidley, T. M., Boardman, J. D., Domingue, B. W. & Boardman, J. D.

- 1 Changing Polygenic Penetrance on Phenotypes in the 20th Century Among Adults in
2 the US Population. *Sci. Rep.* **6**, 30348 (2016).
- 3 37. Walter, S. *et al.* Association of a Genetic Risk Score With Body Mass Index Across
4 Different Birth Cohorts. *JAMA* **316**, 63 (2016).
- 5 38. de Vlaming, R. *et al.* Meta-GWAS Accuracy and Power (MetaGAP) Calculator Shows
6 that Hiding Heritability Is Partially Due to Imperfect Genetic Correlations across
7 Studies. *PLOS Genet.* **13**, e1006495 (2017).
- 8 39. Mehta, D. *et al.* Evidence for Genetic Overlap Between Schizophrenia and Age at First
9 Birth in Women. *JAMA Psychiatry* **73**, 497–505 (2016).
- 10 40. Briley, D. A., Tropf, F. C. & Mills, M. C. What Explains the Heritability of Completed
11 Fertility? Evidence from Two Large Twin Studies. *Behav. Genet.* **47**, 36–51 (2017).
- 12 41. Tropf, F. C. & Mandemakers, J. J. Is the Association Between Education and Fertility
13 Postponement Causal? The Role of Family Background Factors. *Demography* **54**, 71–
14 91 (2017).
- 15 42. Mills, M. C. *Introducing survival and event history analysis.* (Sage Publications,
16 2011).
- 17 43. Tropf, F. C., Barban, N., Mills, M. C., Snieder, H. & Mandemakers, J. J. Genetic
18 influence on age at first birth of female twins born in the UK, 1919-68. *Popul. Stud.*
19 (NY). 129–145. (2015).
- 20 44. Stearns, S. C., Byars, S. G., Govindaraju, D. R. & Ewbank, D. Measuring selection in
21 contemporary human populations. *Nat. Rev. Genet.* **11**, 611–622 (2010).
- 22 45. Courtiol, A., Tropf, F. C. & Mills, M. C. When genes and environment disagree:
23 Making sense of trends in recent human evolution. *Proc Natl Acad Sci US A* **113**,
24 7693–7695 (2016).
- 25 46. Visscher, P. M., Hill, W. G. & Wray, N. R. Heritability in the genomics era—concepts
26 and misconceptions. *Nat. Rev. Genet.* **9**, 255–266 (2008).
- 27 47. Witte, J. S., Visscher, P. M. & Wray, N. R. The contribution of genetic variants to
28 disease depends on the ruler. *Nat. Rev. Genet.* **15**, 765–776 (2014).
- 29 48. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide
30 complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 31 49. Wray, N. R. & Maier, R. Genetic Basis of Complex Genetic Disease: The Contribution
32 of Disease Heterogeneity to Missing Heritability. *Curr. Epidemiol. Reports* **1**, 220–227
33 (2014).

34

35

1 **Acknowledgements:** Funding was provided by grants awarded to M.C.M.: ERC
2 Consolidator Grant SOCIOGENOME (615603), UK ESRC/NCRM SOCGEN grant
3 (ES/N011856/1) and the Wellcome Trust ISSF and John Fell Fund. The funders had no role
4 in study design, data collection and analysis, decision to publish, or preparation of the
5 manuscript.

6

7 **Authors Contributions:**

8 FCT, SHL, MCM developed the study concept and study design

9 MRR, FCT developed the concept for and performed the simulation studies

10 FCT, RMV, GS, CR performed data analysis and visualization

11 TE, AM, SEM, NGM, AN, SHL, AB provided data, input on data analysis and interpretation,
12 and imputed data

13 FCT, MCM, SHL, MRR, HS, GS, RV drafted the manuscript

14 FCT, MCM, MRR, RMV, GS, PJM, RV, NGM, NB, DAB, CR, RH revised the manuscript

15 All authors approved the final version of the manuscript for submission.

16 **Competing Interests:** The authors declare no competing interests.

17

Boxes

Box 1: Definitions of heritability

Heritability

Heritability is the proportion of the phenotypic variance accounted for by genetic effects and narrow sense heritability refers to the additive genetic variance component (for discussion also see ^{6,46}). There are several ways to estimate heritability. First, the highest and prominent estimates are derived from family-based studies (h^2_{family}), such as twin studies, where, typically, the genetic resemblance between relatives is mapped to phenotypic similarity, taking unique- and shared-environment effects into account. Under several assumptions, estimates of h^2_{family} ought to reflect only additive-genetic effects. A second method is the proportion accounted for by genetic variants known to be robustly associated with the phenotype of interest, derived from a GWAS (genome-wide association study) (h^2_{GWAS}). This measure tends to produce the lowest levels. Finally, there is the proportion of phenotypic variance jointly accounted for by all variants on standard GWAS chips. This is sometimes referred to as the SNP- or chip-based heritability (h^2_{SNP}). Typically, h^2_{SNP} is substantially larger than h^2_{GWAS} and provides an ‘upper level estimate’ of the genetic effects that could be identified with a well-powered GWAS. The h^2_{GWAS} increases in tandem with GWAS sample sizes and is expected to approach h^2_{SNP} asymptotically under the assumption that the phenotype of interest is homogeneous in its genetic architecture across different environments.

Missing heritability

The gap between the h^2_{family} and h^2_{GWAS} is referred to as ‘missing heritability’.¹² Potential reasons for missing heritability are for example non-additive genetic effects (although empirical evidence on this is scarce),^{2,13} large effects of rare variants,¹⁵ and potentially inflated estimates from twin studies due to shared environmental factors.¹⁶ The missing heritability is commonly defined as the sum of the still-missing and hidden heritability, which we define below.⁴⁷

Still-missing heritability

Yang and colleagues⁸ argued that most genetic effects are too small to be reliably detected in GWAS of current sample sizes which is why they proposed the whole-genome restricted maximum likelihood estimation performed by GCTA software.⁴⁸ Studies applying these whole-genome methods typically produce estimates that lie between twin studies and polygenic scores $h^2_{\text{GWAS}} < h^2_{\text{SNP}} < h^2_{\text{family}}$. The discrepancy $h^2_{\text{SNP}} < h^2_{\text{family}}$ has been referred to as ‘still-missing heritability’.¹¹ A stylized fact is that for many traits the still-missing heritability is roughly equal to h^2_{SNP} .⁴⁹ It is generally assumed that by genotyping rarer and structural variants, the still-missing heritability will decrease, as the denser arrays will increase h^2_{SNP} .

Hidden heritability

Since we expect to be able to almost fully capture h^2_{SNP} in the long run, the discrepancy between h^2_{SNP} and h^2_{GWAS} is sometimes referred to as ‘hidden heritability’.⁴⁷ The current study is mainly interested in the question how h^2_{SNP} changes, depending on whether we examine differences within or between populations. Here we focus on hidden heritability as the genetic variation due to heterogeneity that cannot possibly be explained by SNP associations based on meta-analyses of multiple populations. Since h^2_{GWAS} is usually inferred from meta-analyses that include multiple populations, heterogeneity in genetic effects on a phenotype between

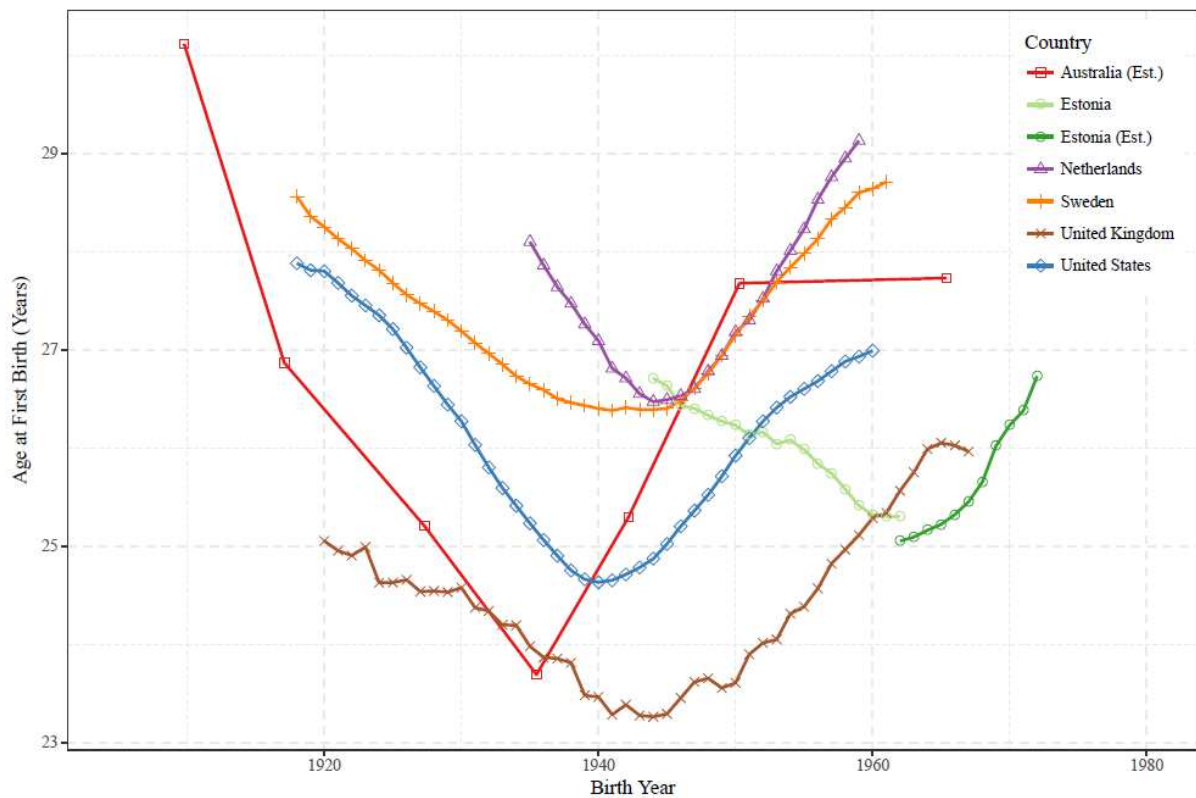
1 these populations could deflate h^2_{GWAS} and would also deflate h^2_{SNP} – which is typically
2 obtained within single populations. Within a single design we therefore demonstrate how one
3 estimate of h^2 depends upon population heterogeneity. Missing heritability is thus commonly
4 defined as the sum of the still-missing and hidden heritability.⁴⁷ As indicated, the hidden
5 portion will decrease as sample sizes grow and the still-missing portion will decrease with
6 denser forms of genotyping.

7

1
2
3

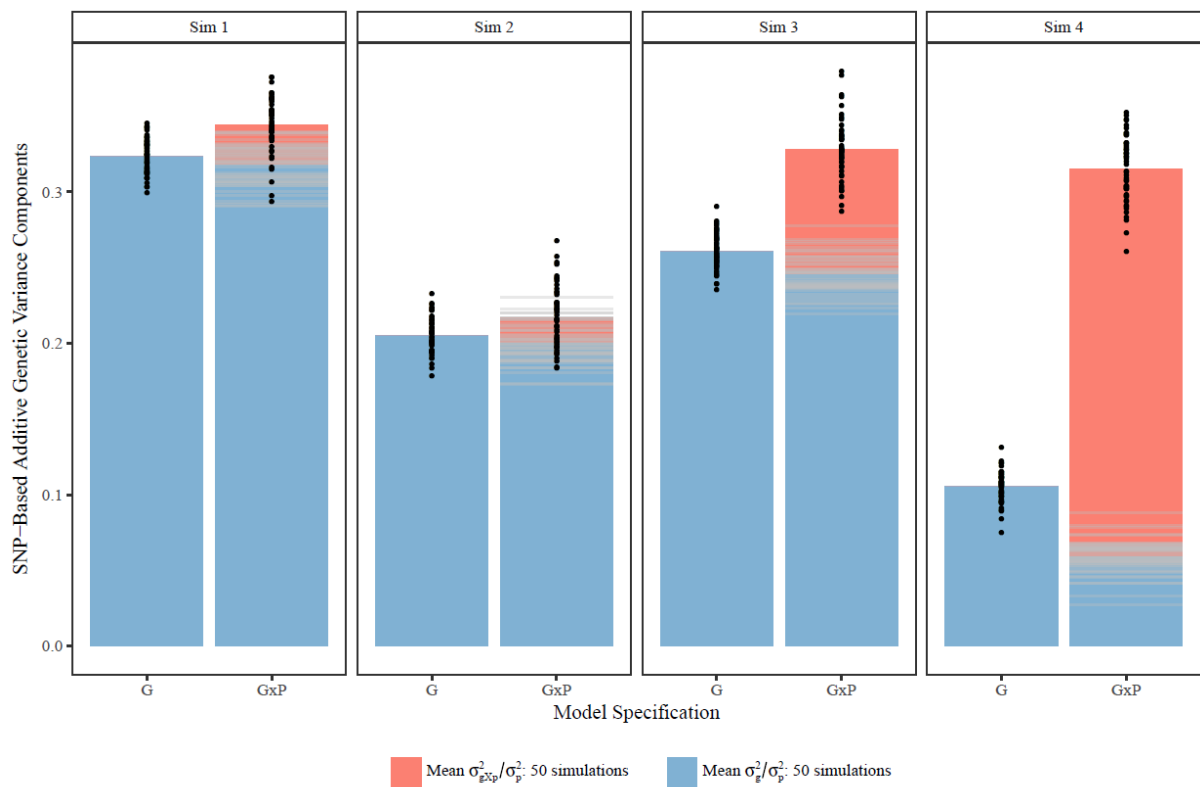
Figure legends:

4 **Figure 1. Trends in mean age at first birth of women indicating environmental changes**
5 **across cohorts (1903-1970) from the US, UK, Sweden, the Netherlands, Estonia and**
6 **Australia.** Trends in the mean age at first birth of women are based on aggregated data
7 obtained from Human Fertility Database and the Human Fertility Collection (for details see
8 Text S3). For Estonia, from 1962 onwards, we used estimated age at first births based on
9 women older than 40. For Australia, no official data was available and the trends have been
10 estimated from the QIMR dataset, averaged for each decade.



11
12

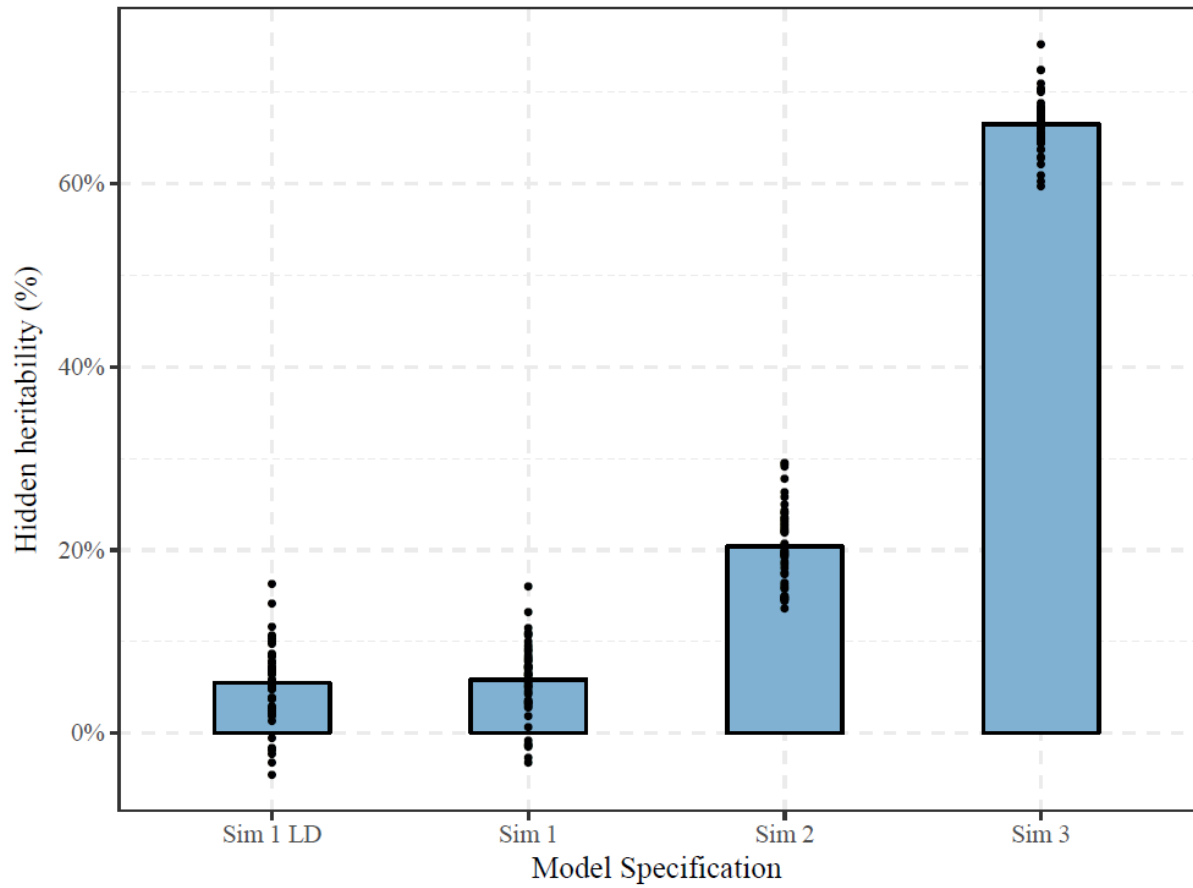
1 **Figure 2. Stacked Bar Charts of average between (σ_g^2) and within (σ_{gxp}^2) variance**
 2 **explanation by common SNPs estimated across 50 simulated phenotypes in two**
 3 **model specifications (standard GREML model and the gene-environment**
 4 **interaction model by study population (GxP) and for four simulated phenotypes:**
 5 **Sim 1 with homogeneous SNP-based heritability 0.5 without gene-environment**
 6 **interaction, Sim 2 heterogeneous SNP-based heritability between 0.25-0.625**
 7 **without gene environment interaction, Sim 3 with homogeneous SNP-based**
 8 **heritability 0.5 with gene-environment interaction (genetic correlation of 0.8**
 9 **across populations) and Sim 4 with homogeneous SNP-based heritability 0.5 with**
 10 **gene-environment interaction (genetic correlation of 0.5 across populations).**
 11 **Individual model estimates are represented by black dots, individual σ_g^2**
 12 **components in the GxP models in gray stripes.**



13

14

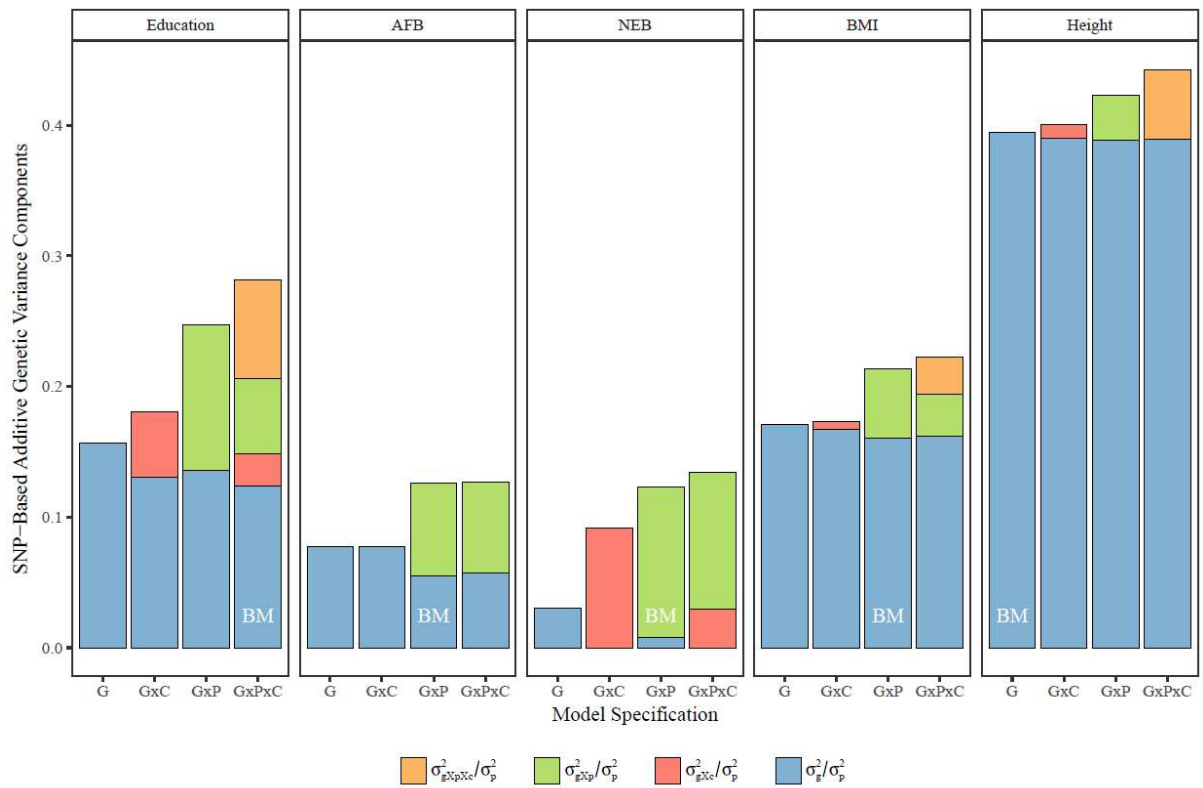
1 **Figure 3. Bar Charts of average % of hidden heritability due to heterogeneity (% of**
2 **h^2_{SNP} of the best fitting model which is not captured in standard GREML models)**
3 **for Sim 1 including and excluding causal variants (Sim LD), for Sim 3 and 4.**
4 **Individual estimates are represented by black dots.**



5

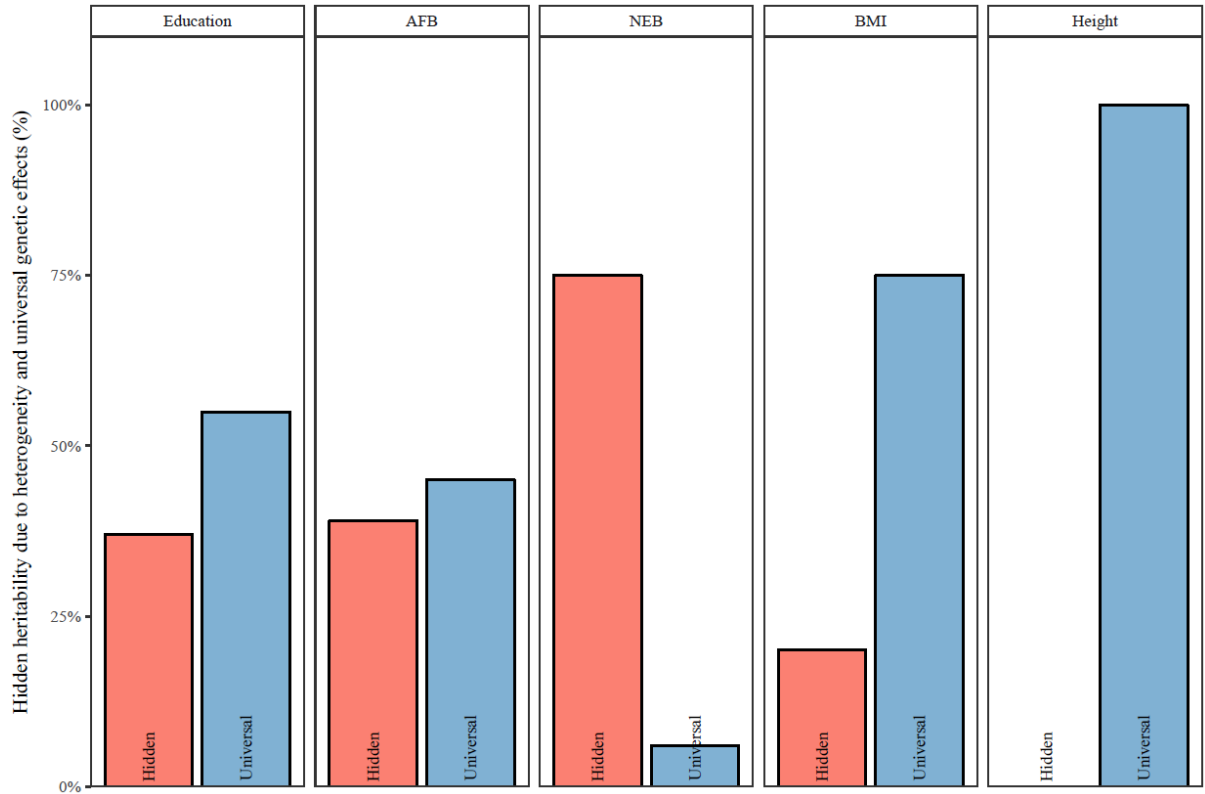
6

1 **Figure 4. Stacked Bar Charts of average between (σ_g^2) and within ($\sigma_{gXp}^2, \sigma_{gXc}^2, \sigma_{gXpXc}^2$)**
2 **variance explanation by common SNPs estimated for Height, BMI, education, age**
3 **at first birth (AFB) and number of children (NEB) in four model specifications**
4 **(G, GxP, GxC, GxPxC).** The best model (BM in white, in chart) is based on
5 likelihood ratio tests comparing the full model with one constraining the respective
6 variance component to 0; see Table S2. σ_g^2/σ_p^2 = proportion of observed variance in the
7 outcome associated with genetic variance across all environments, $\sigma_{gXp}^2/\sigma_p^2$ =
8 proportion of observed variance in the outcomes associated with *additional* genetic
9 variance within populations, $\sigma_{gXc}^2/\sigma_p^2$ = proportion of observed variance associated
10 with *additional* genetic variance within demographic birth cohorts, $\sigma_{gXpXc}^2/\sigma_p^2$ =
11 proportion of observed variance associated with *additional* genetic variance within
12 populations and demographic birth cohorts. Models specifications G, GxP, GxC,
13 GxPxC refer to the model specifications including the respective variance components
14 as well as those of lower order – see Material and Methods. For detailed results see
15 Table S3-7.



16
17

1 **Figure 5. Bar Charts of average % of hidden heritability due to heterogeneity (% of**
 2 **h^2_{SNP} of the best fitting model which is not captured in standard GREML models)**
 3 **and of universal genetic effects (% of h^2_{SNP} of the best fitting model which is**
 4 **effectively identical across the defined environments)**



5
6