

# Imprint of Assortative Mating on the Human Genome

Loic Yengo<sup>1,\*</sup>, Matthew R. Robinson<sup>1,2</sup>, Matthew C. Keller<sup>3</sup>, Kathryn E. Kemper<sup>1</sup>, Yuanhao Yang<sup>1</sup>, Maciej Trzaskowski<sup>1</sup>, Jacob Gratten<sup>1</sup>, Patrick Turley<sup>4,5</sup>, David Cesarini<sup>6,7,8</sup>, Daniel J. Benjamin<sup>9,6,10</sup>, Naomi R. Wray<sup>1,11</sup>, Michael E. Goddard<sup>12,13</sup>, Jian Yang<sup>1,11</sup> & Peter M. Visscher<sup>1,11,\*</sup>

<sup>1</sup>Institute for Molecular Bioscience, The University of Queensland, Brisbane 4072, Australia; <sup>2</sup>Department of Computational Biology, University of Lausanne, CH-1015, Switzerland. <sup>3</sup>Department of Psychology & Neuroscience, Institute for Behavioral Genetics, University of Colorado at Boulder. <sup>4</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts, USA. <sup>5</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA. <sup>6</sup>National Bureau of Economic Research, Cambridge, MA, USA. <sup>7</sup>Department of Economics, New York University, New York, New York, USA. <sup>8</sup>Center for Experimental Social Science, New York University, New York, New York, USA. <sup>9</sup>Center for Economic and Social Research, University of Southern California, Los Angeles, California, USA. <sup>10</sup>Department of Economics, University of Southern California, Los Angeles, California, USA. <sup>11</sup>Queensland Brain Institute, The University of Queensland, Brisbane 4072, Australia; <sup>12</sup>Faculty of Veterinary and Agricultural Science, University of Melbourne, Parkville, Victoria, Australia; <sup>13</sup>Biosciences Research Division, Department of Economic Development, Jobs, Transport and Resources, Bundoora, Victoria, Australia.

\*To whom correspondence may be addressed. Email: [l.yengodimbou@uq.edu.au](mailto:l.yengodimbou@uq.edu.au) or

[peter.visscher@uq.edu.au](mailto:peter.visscher@uq.edu.au)

**Key words:** assortative mating; mate choice; gametic phase disequilibrium; quantitative Genetics; Single Nucleotide Polymorphism; summary statistics.

**Non-random mate-choice with respect to complex traits is widely observed in humans, but whether this reflects true phenotypic assortment, environment (social homogamy) or convergence after choosing a partner is not known. Understanding the causes of mate choice is important, because assortative mating (AM) if based upon heritable traits, has genetic and evolutionary consequences. AM is predicted under Fisher's classical theory<sup>1</sup> to induce a signature in the genome at trait-associated loci that can be detected and quantified. Here, we develop and apply a method to quantify AM on a specific trait by estimating the correlation ( $\theta$ ) between genetic predictors of the trait from SNPs on odd versus even chromosomes. We show by theory and simulation that the effect of AM can be distinguished from population stratification. We applied this approach to 32 complex traits and diseases using SNP data from ~400,000 unrelated individuals of European ancestry. We found significant evidence of AM for height ( $\theta=3.2\%$ ) and educational attainment ( $\theta=2.7\%$ ), both consistent with theoretical predictions. Overall, our results imply that AM involves multiple traits, affects the genomic architecture of loci that are associated with these traits and that the consequence of mate choice can be detected from a random sample of genomes.**

Non-random mating in natural populations has short and long-term evolutionary consequences. In many species, including humans, mate choice is often associated with phenotypic similarities between mates<sup>2,3</sup>. Such phenotypic similarities have multiple sources, for example social homogamy, the preference for a mate from the same environment, or because of primary assortment on certain traits observable at the time of mate choice. Contrary to social homogamy, primary phenotypic assortment, here referred to as assortative mating (AM), has genetic and evolutionary consequences and therefore is the focus of our study. In humans, AM involves multiple complex traits<sup>4-8</sup> and can sometimes lead to similar susceptibility to diseases<sup>9-12</sup>. The genetic effects of AM were first studied in the seminal articles of Fisher (1918)<sup>1</sup> and Wright (1921)<sup>13</sup>. Those two founding contributions, further complemented by Crow & Kimura (1970)<sup>14</sup> and others<sup>15-17</sup> have set the basis of the theory of AM on complex traits. AM theory predicts three main genetic consequences of a positive correlation between the phenotypes of mates in a population: (i) an increase of the genetic variance in the population, (ii) an increase in the correlation between relatives and (iii) an increase of homozygosity (deviation from Hardy-Weinberg Equilibrium; HWE), in particular at causal loci. These seemingly distinct consequences are nonetheless explained by the same cause: directional correlation between trait-increasing alleles (TIA), also referred to as gametic phase disequilibrium (GPD), induced both within and between loci (**Fig. 1**). AM-induced GPD implies correlations between physically distant loci (between chromosomes for example) and is thus distinct from linkage disequilibrium. Therefore, AM leads to a genomic signature of trait-associated loci that can be quantified by estimating GPD.

Previous studies<sup>18–20</sup> have been successful at detecting GPD by direct quantification of increased homozygosity at ancestry-associated loci. Beyond ancestry-related traits, such endeavour is particularly challenging for polygenic traits as theory<sup>14</sup> predicts an increase of homozygosity due to AM inversely proportional to the number  $M$  of causal variants<sup>14,21</sup>. For a highly polygenic trait like height with an estimated  $M \sim 100,000$  for common variants<sup>22</sup>, the expected increase in homozygosity would be of the order of  $\sim 1/2M = 5 \times 10^{-6}$ , i.e. negligible (**Supplementary Note 1**). Extremely large studies would therefore be required to quantify systematic deviation from HWE at height-associated single-nucleotide polymorphism (SNP) as shown in a recent study<sup>18</sup> that failed to detect such an effect. Another study<sup>23</sup> in  $\sim 6,800$  participants of European ancestry, reported evidence of deviation from HWE at height associated loci. This study however did not account for within-sample population stratification and therefore their reported estimates are likely inflated for this reason. Overall, study designs using deviation from HWE for quantifying GPD can be successful for detecting ancestry-based AM (ancestral endogamy) because the number of loci distinguishing ancestries is relatively small<sup>24</sup>, and ancestral endogamy is strong<sup>18</sup>, but are less powerful to detect trait-specific AM. In contrast to HWE-based estimation strategies, quantifying GPD on the basis of pairwise correlations between TIAs is much more tractable as the number of pairs of loci involved, of the order of  $\sim M^2$ , compensates for the magnitude of the expected covariance for each pair,  $\sim 1/2M$ . This compensation explains in essence why AM increases the genetic variance in a population<sup>14,21</sup>.

GPD due to AM causes individuals that carry TIAs at one locus to be more likely to carry TIAs at other loci than expected under gametic phase equilibrium.

Consequently, individuals with many TIAs on even chromosomes are likely to have more than average TIAs on odd chromosomes. We quantify this effect by calculating genetic predictors for a trait from the SNPs on odd chromosomes and from the SNPs on even chromosomes and then calculating the correlation ( $\theta$ ) between these two predictors. To calculate these predictors we use estimates of the effect of each SNP on a trait from publicly available summary statistics from genome-wide association studies (GWAS) of large sample size. We apply these estimated SNP effects to individuals in a separate sample who have SNP genotypes available. We can calculate the trait predictor based on odd and even chromosomes separately and estimate the correlation between them (i.e.  $\theta$ ). Our method measures the effect on the genome due to AM in previous generations and thus does not require observed phenotypes or the use of mate pairs. Under the null hypothesis of random mating (RM), the correlation between alleles on different chromosomes is expected to be 0 as a consequence of the independent segregation of chromosomes. However, population stratification can induce spurious correlations between alleles, even at distant loci. Intuitively, if  $\theta$  is estimated from a mixture of two sub-populations with distinct allele frequencies, then having TIAs more frequent in one of the sub-populations (even by chance) would result in an apparent correlation between TIAs even when such correlation is absent in each sub-population (**Supplementary Note 2**). We show through simulations how the effect of population stratification can be corrected with our method. We then applied our method to estimate AM-induced GPD for 32 traits and diseases in samples of unrelated genomes from three independent cohorts: ~350,000 participants of the UK Biobank (UKB), ~54,000 participants of the Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort and ~8,500 participants of the US Health and

Retirement Study (HRS). We find evidence of AM for a number of complex traits, including height and educational attainment.

## Results

### *Theory underlying the estimation of GPD from SNP data*

We derived (**Supplementary Note 1**) the expected value of the correlation across individuals between the trait predictors from SNPs on odd ( $S_o$ ) and even ( $S_e$ ) chromosomes as a function of the phenotypic correlation between mates ( $r$ ), the equilibrium heritability of the trait ( $h_{eq}^2$ ), the fraction of that heritability captured by SNPs ( $f_{eq}$ ), the sample size ( $N$ ) of the reference GWAS in which effect sizes are estimated using classical linear regression, one SNP at a time; and the number of causal loci ( $M$ ) contributing to the trait (which differs from the number of statistically associated loci). The main result is that for a large number of trait loci,

$$\theta = \frac{\rho f_0}{2 - \rho(2 - f_0)} \left[ 1 + \frac{M(1 - \rho)}{N h_{eq}^2} \left( 1 + \frac{\rho f_0}{2(1 - \rho)} \right)^{-3} \right]^{-1} \quad (1)$$

with  $\rho = r h_{eq}^2$  being the correlation between additive genetic values of mates expected under AM<sup>17</sup> and  $f_0 \approx f_{eq}/(1 - \rho)$ , the fraction of heritability captured by SNPs in the base population (**Supplementary Note 1**). These parameters do not need to be known or estimated, but can be used to provide *a priori* expectation of  $\theta$  or *a posteriori* rationalisation. Hence, quantification of GPD can be directly obtained from estimates of  $\theta$  using empirical data. For the sake of simplicity, we hereafter refer to estimates of  $\theta$  as estimates of GPD. Equation (1) implies that the expected correlation

$\theta$  between  $S_o$  and  $S_e$  increases with  $N$ , i.e. with better estimation of SNP effects from the reference GWAS, and with  $f_{eq}$ , i.e. with better tagging of causal variants underlying the full narrow sense heritability.

We derived (**Supplementary Note 2**) that estimates of  $\theta$  from the regression of  $S_o$  on  $S_e$  can be inflated by population stratification, especially when TIAs are highly differentiated between sub-populations. We performed a number of simulations (**Supplementary Notes 2 and 3**) to validate the impact of population stratification on our estimator of GPD and show how to adjust for it using principal components derived separately for odd and even chromosomes (**Fig. S1** and **Fig. S2**). We used this strategy to quantify GPD in real data, and therefore adjusted all GPD estimates for 20 genotypic principal components to correct within sample population stratification (**Materials and Methods**). Also, given that most GPD estimates are small, all GPD estimates (correlations) reported below are expressed as percentages (e.g. 3% instead of 0.03).

### ***Quantifying AM-induced GPD in complex traits***

We first analysed height and educational attainment (EA), two reference traits with long-standing evidence of a positive correlation between mates. For height, we used estimated effect sizes, from summary statistics of the latest GWAS meta-analysis of the GIANT consortium (Wood et al. 2014)<sup>25</sup>, of 9,447 near independent HapMap3 (HM3) SNPs selected using the LD clumping algorithm implemented in the software PLINK<sup>26</sup> (linkage disequilibrium (LD)  $r^2 < 0.1$  for SNPs  $< 1$  Mb apart and association  $p$ -value  $< 0.005$ ). We thus selected these SNPs to be enriched for true association with height. We estimated in UKB participants the correlation between height increasing

alleles on odd versus even chromosomes to be  $\theta_{\text{height}}=3.0\%$  (s.e. 0.2%; **Fig. 2**) and replicated this finding in GERA ( $\theta_{\text{height}}=4.1\%$ , s.e. 0.4%; **Fig. 2**) and HRS ( $\theta_{\text{height}}=4.4\%$ , s.e. 1.1%; **Fig. 2**). A meta-analysis of these three estimates yields a combined GPD among height-increasing alleles of 3.2% (s.e. 0.2%,  $p=6.5\times 10^{-89}$ ). To dismiss possible biases due to cryptic sample overlap or residual population stratification in summary statistics from the Wood *et al.* study, we re-estimated  $\theta_{\text{height}}$  using summary statistics of a family-based GWAS that provide stringent control for population stratification<sup>27</sup>. We therefore meta-analysed summary statistics from the Robinson *et al.* (2015) study<sup>27</sup> in 17,500 quasi-independent sib pairs with that from a similar analysis performed in 21,783 quasi-independent sib-pairs identified in the UKB (**Materials and Methods**). Using effect sizes of the 9,447 selected SNPs, re-estimated in the combined family-based GWAS, we found consistent GPD estimates between UKB not including sib-pairs ( $\theta_{\text{height}}=2.1\%$ , s.e. 0.2%;  $p=8.4\times 10^{-36}$ ), GERA ( $\theta_{\text{height}}=2.1\%$ , s.e. 0.4%;  $p=1.4\times 10^{-6}$ ) and HRS ( $\theta_{\text{height}}=2.5\%$ , s.e. 1.1%;  $p=0.02$ ). The meta-analysis of these three estimates yields  $\theta_{\text{height}}=2.1\%$  (s.e. 0.2%;  $p=4.7\times 10^{-42}$ ). Note that lower estimates (2.1% vs 3.4%) are expected here because of the smaller sample size ( $N=39,283$ ) of this family-based GWAS, as predicted by equation (1).

For EA, we used estimated effect sizes, from the summary statistics of a large GWAS meta-analysis of the number of years of education (Okbay *et al.* 2016)<sup>28</sup>, of 4,618 near independent HM3 SNPs selected using the same LD clumping procedure described above. Using genotypes of 238,193 UKB participants not included in the Okbay *et al.* study (**Materials and Methods**), we found that  $\theta_{\text{EA}}=2.9\%$  (s.e. 0.2%; **Fig. 2**) and replicated this finding in GERA ( $\theta_{\text{EA}}=1.8\%$ , s.e. 0.4%; **Fig. 2**). We also attempted replication in HRS but the estimate we found ( $\theta_{\text{EA}}=8.9\%$ , s.e. 1.1%; **Fig. 2**) was likely



inflated given that HRS was part of the Okbay *et al.* meta-analysis (**Supplementary note 4**). We therefore only meta-analysed GPD estimates from UKB and GERA and found the average correlation between EA increasing alleles on odd versus even chromosomes to be  $\theta_{EA}=2.7\%$  (s.e. 0.3%,  $p=6\times 10^{-46}$ ; **Fig. 2**). We also re-estimated the effect sizes of the 4,618 selected SNPs on EA, using the same within-family procedure described above. We found GPD estimates of  $\sim 0.4\%$  (s.e. 0.4%) in GERA and  $\sim 0.3\%$  (s.e. 0.1%) in UKB participants unrelated with any of the 21,783 sib-pairs used to estimate effect sizes. The meta-analysis of the latter two estimates is  $\theta_{EA}=0.31\%$  (s.e. 0.16%;  $p=0.05$ ). As shown below, this lower estimate is expected as the consequence of the smaller sample size used to estimate SNPs effects.

We compared GPD estimates with theoretical predictions of  $\theta$  from equation (1). Equation (1) predicts  $\theta$  from the sample size of the reference GWAS ( $N=253,288$  for height and 293,723 for EA), the correlation between mates, the equilibrium heritability (80% and 40% for height and EA respectively<sup>29</sup>), the number of causal variants SNPs ( $M\sim 100,000$  assumed for both traits) and the heritability captured by SNPs used to estimate  $\theta$ . Using  $\sim 1,000$  unrelated trios (two parents and one offspring) from UKB<sup>30</sup> we estimated the correlation between mates for height and EA to be 0.24 (s.e. 0.03) and 0.35 (s.e. 0.03), respectively. We estimated the SNP heritability captured by each set of SNPs used to estimate  $\theta$  in HRS using the software GCTA<sup>31</sup>, resulting in estimates of  $h^2_{\text{height}} = 27.3\%$  (s.e. 1.7%) and  $h^2_{EA} = 15.1\%$  (s.e. 1.3%). With these five input parameters, equation (1) predicts  $\theta$  to be  $\sim 3.2\%$  for height and  $\sim 1.9\%$  for EA. Using estimated effective sample sizes of within-family GWAS ( $N_{\text{eff}} = 39,283$  for height and 15,559 for EA; **Materials and Methods**), equation (1) predicts  $\theta$  to be  $\sim 1.3\%$  for height and  $\sim 0.2\%$  for EA. Our estimates of GPD among trait-

associated ( $\theta_{\text{height}}=3.2\%$ , s.e. 0.2;  $\theta_{\text{EA}}=2.7\%$ , s.e. 0.3%) are therefore consistent with these predictions. Everything held constant, equation (1) also predicts that with a much larger sample size of the discovery GWAS, for instance  $>1,000,000$  participants,  $\theta_{\text{height}}$  would be  $\sim 4\%$  and  $\theta_{\text{EA}} \sim 3\%$ .

We extended our primary analysis of height and EA to detect GPD in 30 additional complex traits and diseases (**Table S1**) using the same strategy. Among these traits, we analysed bone mineral density (BMD)<sup>32</sup> as a null trait given that non-significant mate correlation was previously reported<sup>33</sup>. As expected, we did not find significant GPD associated with BMD (meta-analysis of UKB, GERA and HRS:  $\theta_{\text{BMD}}=0.09\%$ , s.e. 0.2%;  $p=0.64$ ). After Bonferroni correction applied to the meta-analysis of UKB, GERA and HRS ( $p<0.05/32 \approx 1.56 \times 10^{-3}$ ), we did not detect significant GPD for any of these other traits. We believe that this observation is likely explained by lack of statistical power, in particular resulting from the smaller sizes of GWAS used for these traits (on average  $\sim 73,000$  participants compared to  $\sim 273,000$  on average for height and EA) or from smaller variance explained by SNPs selected to calculate genetic predictors of these traits. As an example, although the GWAS of body mass index (BMI) used in this study is similar in size with that of height (**Table S3**), our estimation in HRS participants of the variance explained by the 2,362 SNPs BMI-associated SNPs selected (**Table S1**) is only  $\sim 6.2\%$  (s.e. 0.9%) versus  $\sim 27.3\%$  (s.e. 1.7%) for height. A much larger GWAS would therefore be required to detect any GPD among BMI-associated alleles using our method.

### *Confirmation using mate pairs*

Another experimental design to quantify the genetic effect of AM on a particular trait consists of estimating the correlation of genetic predictors of this trait between mates<sup>33–35</sup>. We quantified the mate correlation ( $r_m$ ) of genetic predictors of all 32 traits (**Table S2**) using 18,984 unrelated couples identified in the UKB (**Materials and Methods**). We found significant correlations between mates for genetic predictors of height ( $r_m=5.9\%$ , s.e. 0.8%,  $p=9.2\times 10^{-14}$ ) and EA ( $r_m=6.1\%$ , s.e. 0.9%,  $p=7.3\times 10^{-11}$ ). Across all traits, we estimated the regression slope of  $r_m$  estimates onto  $\theta$  estimates to be 1.8 (s.e. 0.2) (**Fig. 3**). Both these results are consistent with our derivation that the expected mate correlation of genetic predictors is approximately twice the expected value of  $\theta$  (**Supplementary note 4**).

## Discussion

We have shown in this study that the genomic signature of AM can be detected and quantified using SNP data in a random sample of genomes from the population, even in the absence of data on mate pairs. This is an important aspect of our method since large datasets on mate pairs are rare and may be absent in natural populations. We confirm the genetic basis for AM for height and EA, consistent with the assumption of primary assortment on these traits. We showed that our estimates of GPD from real data are consistent with theoretical predictions made under simplifying assumptions such as equal SNP effect sizes, population being at equilibrium and that the number of common causal variants of the order of  $\sim 100,000$  (**Supplementary Note 1**). We did not however detect significant GPD for the other traits and diseases analysed in this study. Beyond true negatives such as bone mineral density, we believe that the relatively smaller size of GWAS used in our inference has reduced the power to

detect the genetic signature of AM in more traits and diseases. We cannot therefore draw firm conclusion from our observations on the importance of primary assortment (as opposed to environmental correlation) to the resemblance between mates for some of these traits such as smoking habits<sup>36</sup>, alcohol consumption<sup>36</sup> or susceptibility to psychiatric disorders<sup>12</sup>. Overall, our methodology is straightforward and can be applied to a wider variety of traits and in other species, provided that summaries of trait-variant associations are available. AM is multi-dimensional in essence as mate choice depends on multiple physical and behavioural traits which may or may not share a common genetic basis<sup>5,37</sup>. Studying networks of traits and genes driving AM is one of the challenges to meet for improving our understanding of the genetic consequences of AM in a population. As a step in this direction, our method can be for example applied to quantify consequences of AM on gene expression or at any other molecular level, through the use of SNP predictors of these endogenous traits.

Our study predicts that for traits affected by AM the estimates of SNP effects from within-family experimental designs tend to be smaller than those from a population sample, and that a genetic predictor generated from a population sample will explain less variation than expected when applied to a population not undergoing assortative mating.

Our study has a number of limitations. The first one is that certain aspects of our approach are very conservative. We have attempted to quantify GPD induced by AM while applying stringent correction for population stratification. Although such a strategy is expected to yield unbiased estimates it still faces the difficulty of distinguishing AM on population stratification related features from AM on trait

specific features. Height is a typical example. AM on height occurs but, in addition, people tend to marry within geographical sub-populations (countries for example) which differ in mean height<sup>27</sup>. Correcting for population stratification using principal components would consequently remove part of the signal that we want to detect. We have nevertheless been able to detect GPD among height increasing alleles as a consequence of the large sample size of the discovery GWAS, the strength of assortment of mates, and the high heritability of this trait.

The second limitation relates to our strategy for SNP selection. We have included in our analyses SNPs using a low and arbitrary threshold ( $p < 0.005$ ) on the significance of association with the trait. Although this strategy is not expected to bias the covariance between  $S_e$  and  $S_o$ , it may increase both their variances and thus potentially induces downward biases in GPD estimates. We chose nonetheless this strategy to maximize the fraction of heritability captured by SNPs, which influences the expected correlation between  $S_e$  and  $S_o$  as derived in equation (1). As an example, if GPD is inferred using genome-wide significant SNPs from Okbay *et al.* (2016), which explain ~3% of the variance of EA, the expected correlation between  $S_e$  and  $S_o$  would only be ~0.45% under assumptions made above. Such small correlation is nearly undetectable in cohorts with less than 300,000 participants (**Materials and Methods**). Another SNP selection strategy could have been used to reach a better trade-off between bias and power but this would generally require observed phenotypes to optimize genetic predictors<sup>33,34</sup> (find the best  $p$ -value threshold or shrinkage parameter).

In conclusion, our study provides empirical quantification of GPD induced among trait-associated alleles, a phenomenon predicted by theory exactly a century ago by Fisher (1918)<sup>1</sup>. The human genome has a pattern of trait-associated loci that is shaped by human behaviour (mate choice), as well as natural selection<sup>33,38–40</sup>. The imprint of assortative mating on the contemporary human genome reflects mate choice in the 1930-1970s and in generations prior to that. Although there is much more mobility within and between human populations in the 21<sup>st</sup> century, the underlying factors that determine mate choice remain stable<sup>11,35</sup>, and we may expect to continue to see its effect in the genome.

## Materials and Methods

### *Estimation of GPD from SNP data*

Our inference of GPD in a given sample of genomes is based on the correlation  $\theta$  between polygenic scores  $S_e$  and  $S_o$  calculated from SNPs on even and odd chromosomes respectively. For each individual from the study population, these scores are obtained as linear combinations of SNPs allele counts weighted by their estimated effect sizes from publicly available GWAS of complex traits and diseases (**Supplementary Note 1**). We used publicly available summary statistics (regression coefficients for each tested SNP and  $p$ -values) from large GWAS on 32 traits (**Table S3**; URLs to download these summary statistics are given in **Supplementary Note 1**). These include GWAS on cognitive traits (educational attainment, intelligence quotient), anthropometric traits (height, body mass index, waist-to-hip ratio), psychiatric disorders (attention deficit hyperactivity disorder, autism spectrum disorder, bipolar disorder, anxiety, major depressive disorder, post-traumatic stress disorder and schizophrenia), other common diseases (coronary artery disease, type 2 diabetes, Crohn's disease and rheumatoid arthritis), blood pressure, reproductive traits, personality traits, alcohol and smoking, and bone mineral density as a null trait. It is important that the sample of people whose genotypes were used was independent of the sample of people used to estimate SNP effects on each trait. Otherwise large biases can be expected as shown in **Supplementary Note 4**. We applied LD score regression (LDSC) for quality control and only kept summary statistics with a corresponding ratio statistic (ratio = (LDSC intercept - 1) / (mean chi-square statistic over all SNPs - 1)) non-significant from 0 (i.e. estimate / standard error < 2) or < 0.2 (**Table S3**). Significance of the GPD estimates was assessed using  $p$ -values from Wald-tests, with the null hypothesis " $H_0: \theta = 0$ " versus the alternative " $H_1: \theta \neq 0$ ".

### ***Correction of population structure***

We used genotypic principal components to correct for population stratification. We calculated 20 principal components from 70,531 near independent SNPs (35,301 on odd chromosomes and 35,230 on even chromosomes) selected using the LD pruning algorithm implemented in PLINK ( $r^2 < 0.1$  for SNPs  $< 1\text{Mb}$  apart). We denote these principal components as PCO for SNPs on odd chromosomes and PCE for SNPs on even chromosomes. When  $\theta$  is estimated from the regression of  $S_e$  onto  $S_o$ , the effect of population stratification is corrected by adjusting the regression for PCOs (and vice versa). This can be summarised using the following regression equations:  $S_e = \theta S_o + \text{PCO}_1 + \dots + \text{PCO}_{20}$  or  $S_o = \theta S_e + \text{PCE}_1 + \dots + \text{PCE}_{20}$ . Since  $S_e$  and  $S_o$  may not have exactly the same variance as a result of SNPs sampling, we chose to estimate  $\theta$  from the regression onto the genetic predictor with the larger variance. We observed nonetheless that estimates obtained from the regression of  $S_e$  onto  $S_o$  are very similar to those obtained from regression of  $S_o$  and  $S_e$  (**Fig. S3**). In the simulation studies (**Supplementary Note 2**) we also considered the case where principal components are calculated from all SNPs available (odd and even chromosomes). In our simulations, principal components were calculated using R version 3.1.2, while in real data principal components were calculated using the fast PCA approach<sup>41</sup> implemented in PLINK version 2.0.

### ***Statistical power***

Theory underlying statistical power to detect correlation is well established<sup>42</sup>. We used equation (2) derived from Ref.<sup>42</sup> to determine the smallest correlation detectable with at least  $1 - \beta = 80\%$  of statistical power at a significance level of  $\alpha = 5\%$ :



$$\log\left(\frac{1+\theta}{1-\theta}\right) = \frac{2|z_{\alpha/2}+z_{\beta}|}{\sqrt{n-3}} \quad (2)$$

In equation (2),  $n$  represents the size of the sample used to estimate  $\theta$ , and  $z_{\alpha/2}$  and  $z_{\beta}$  are the  $\alpha/2$ - and  $\beta$ -upper quantiles of the standard normal distribution (mean 0 and variance 1). With  $\alpha = 5\%$  and  $\beta = 20\%$ ,  $z_{\alpha/2} \approx 1.96$  and  $z_{\beta} \approx 0.84$ . We can therefore detect GPD as small as 1.2% and 0.5% in GERA and UKB respectively, and 0.4% for the meta-analysis of UKB and GERA. For the analysis of mate-pairs we can detect correlation as small as 1.5%.

### ***SNP Genotyping***

#### *UK Biobank data*

We used genotyped and imputed allele counts at 1,312,100 HM3 SNPs in 487,409 participants of the UKBiobank<sup>30,43</sup> (UKB). Extensive description of data can be found here<sup>27</sup>. We restricted the analysis to participants of European ancestry and SNPs with imputation quality  $\geq 0.3$ , minor allele frequency (MAF)  $\geq 1\%$  and Hardy-Weinberg equilibrium test  $p$ -value  $> 10^{-6}$ . Ancestry assignment was performed as follows. We calculated the first two principal components from 2,504 participants of the 1,000 Genomes Project<sup>44</sup> with known ancestries. We then projected UKB participants onto those principal components using SNP loadings of each principal component. We assigned each individual to one of five super-populations in the 1,000 Genomes data: European, African, East Asian, South Asian and Admixed. Our algorithm calculated, for each participant, the probability of membership to the European super-population conditional on their principal components coordinates. The 456,426 out of the original 487,409 participants who had a probability of membership  $> 0.9$  for the European cluster were assigned to the European super-population. Next, to obtain a sample of conventionally unrelated individuals, we estimated the genetic relationship matrix

(GRM) for individuals in the subsample of Europeans using GCTA<sup>31</sup> version 1.9. We iteratively dropped one member from each pair of individuals whose estimated relationship coefficient exceeded 0.05, until no pairs of individuals with a relationship coefficient above 0.05 remained in the sample. This restriction resulted in a sample of 348,502 conventionally unrelated Europeans. In total, we included 348,502 participants in this analysis and 1,124,803 SNPs. The North West Multi-centre Research Ethics Committee (MREC) approved the study and all participants in the UKB study analyzed here provided written informed consent.

#### *Genetic Epidemiology Research on Adult Health and Aging (GERA) cohort data*

We analyzed 60,586 participants of the GERA cohort using genotype data only. Ancestry was assigned using a procedure similar to that described for UKB. Genotype quality control involved standard filters (exclusion of SNPs with missing rate  $\geq 0.02$ , Hardy-Weinberg equilibrium test  $p$ -value  $> 10^{-6}$  or minor allele count  $< 1$ , and removing individuals with missing rate  $\geq 0.02$ ). We imputed genotypes data to the 1,000 Genomes reference panel using IMPUTE2 software. We used GCTA to estimate the GRM of all participants using HM3 SNPs (MAF  $\geq 0.01$  and imputation INFO score  $\geq 0.3$ ). We finally include in the analysis 53,991 unrelated (GRM  $< 0.05$ ) European participants with genotypes at 1,163,290 HM3 SNPs.

#### *Health and Retirement Study (HRS)*

We analysed 8,552 unrelated (GRM  $< 0.05$ ) participants of the HRS cohort. GRM was calculated from 1,118,526 SNPs HM3 imputed to the 1,000 Genomes reference panel using IMPUTE2 software. SNP and samples quality control was similar to what described above for GERA.

### *SNP selection*

We used the LD clumping algorithm implemented in PLINK to identify for each trait near independent SNPs (LD threshold  $r^2 < 0.1$  for SNPs  $< 1$  Mb apart and association  $p$ -value  $< 0.005$ ). LD clumping was performed using genotypes from HRS participants. We restricted the analysis to 1,060,523 HM3 SNPs that passed all quality controls in UKB, GERA and HRS datasets.

### *Sample overlap*

The Okbay *et al.* (2016) GWAS of educational attainment, the Sniekers *et al.* (2017) GWAS of intelligence quotient *al.* (2017) and the Kemp *et al.* (2017) GWAS of bone mineral density, included ~150,000 participants of the UKB (first release of genotypes). To avoid bias due to sample overlap, analyses performed in UKB for these traits were restricted to 238,193 unrelated participants (UKB release 2 only), who also were not related to any of the participants included in those studies (UKB release 1). Participants of the HRS cohort were included in the Okbay *et al.* (2016) study as well as in the Day *et al.* (2015) GWAS of the onset of menopause. For the other GWAS considered in this study, no sample overlap with UKB, GERA or HRS was reported.

### **Sib pairs**

#### *Selection*

We used 21,783 sib pairs of European ancestry previously identified in the UKB<sup>30</sup> using identity-by-descent sharing estimated from SNP data. We applied the within-family QFAM procedure of PLINK, as in Robinson *et al.* (2015)<sup>27</sup>, to assess the

association between HM3 SNPs and height and EA. When applied to sib pairs, this procedure is equivalent to regressing the difference of height or EA between sibs onto the difference of allele counts. These estimates are therefore robust to population stratification. For height, we moreover performed a sample size weighted meta-analysis of estimates from the Robinson *et al.* (2015) study in 17,500 quasi-independent sib pairs, with those obtained in the UKB and used these newly estimated effect sizes to re-estimate GPD in UKB (not including any of the sib-pairs), GERA and HRS. In total we used 21,783 sib-pairs for EA and 39,283 sib-pairs for height.

### *Effective sample size*

We defined the effective sample size ( $N_{eff}$ ) of within-family GWAS using  $N_{pairs}$  independent sib pairs as the sample size of a standard GWAS (where SNP effect are estimated from simple linear regression) such that the estimated SNP effects from the within-family GWAS have the same expected sampling variance as the estimated SNP effects from standard GWAS. This leads to the following equation (derived in

### **Supplementary Note 4)**

$$N_{eff} = N_{pairs} / [2(1 - r_{pairs})] \quad (3)$$

In equation (3),  $r_{pairs}$  represents the phenotypic correlation between sibs. We observed between sibs identified in UKB a correlation  $\sim 0.5$  for height and  $\sim 0.3$  for educational attainment. Therefore, the corresponding effective sample sizes for the within-family GWAS of height and EA are  $39,283 / (2 \times (1 - 0.5)) = 39,283$  and  $21,283 / (2 \times (1 - 0.3)) = 15,559$ .

### **Mate pairs**

We first used 999 unique mate pairs from 1,065 trios composed of both parents and one child, identified among UKB participants using identity-by-descent sharing estimated from SNP data. Details about software and algorithms used to identify those trios are given in Ref.<sup>30</sup> To increase power, we also used household sharing information to identify putative spouse pairs among UKB participants with European ancestry. We therefore selected 18,984 (including 117 from the trios) sex-discordant pairs of participants, recruited from the same centre, who reported living with their spouse or partner in the same type of accommodation, at the same location (east and north coordinates rounded to 1 kilometre), for the same amount of time, with the same number of people in the household, with the same household income, with the same number of smoker in the household, with the same Townsend deprivation index and with a genetic relationship  $< 0.05$ .

### **Data availability**

We used genotypic data from the Resource for Genetic Epidemiology Research on Adult Health and Aging (GERA: dbGaP phs000674.v2.p2), genotypic and phenotypic from the Health and Retirement Study (HRS: dbGaP phs000428.v1.p1), as well as genotypic and phenotypic data from the UK Biobank under project 12505.

### **Acknowledgements**

This research was supported by the Australian Research Council (DP130102666, DP160103860, DP160102400), the Australian National Health and Medical Research Council (1078037, 1078901, 1103418, 1107258, 1127440 and 1113400), the National Institute of Health (NIH grants R01AG042568, P01GM099568 and R01MH100141), the Sylvia & Charles Viertel Charitable Foundation. The Genetic Epidemiology

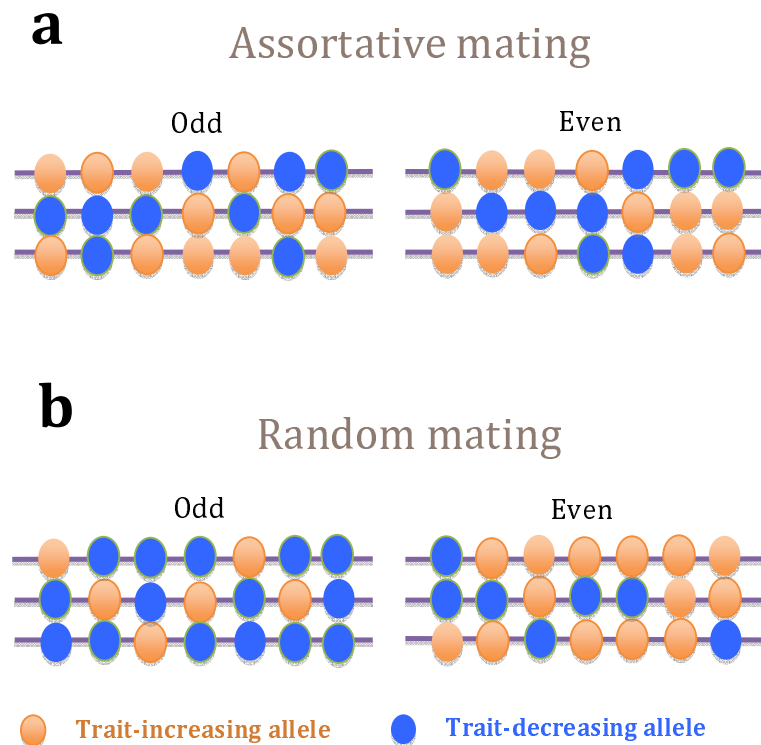
Research on Adult Health and Aging study was supported by grant RC2 AG036607 from the National Institutes of Health, grants from the Robert Wood Johnson Foundation, the Ellison Medical Foundation, the Wayne and Gladys Valley Foundation and Kaiser Permanente. The authors thank the Kaiser Permanente Medical Care Plan, Northern California Region (KPNC) members who have generously agreed to participate in the Kaiser Permanente Research Program on Genes, Environment and Health (RPGEH). This research has been conducted using the UK Biobank Resource under project 12505. We thank Bill Hill for helpful comments and suggestions on the manuscript.

#### **Author contributions**

P.M.V, L.Y., M.R.R, J.Y. and M.E.G. conceived and designed the study. L.Y., M.T. and N.W. curated summary statistics. L.Y. and P.M.V derived the theory. Y.Y., M.T., J.G., K.E.K and L.Y. performed mate pairs analyses. M.C.K., P.T., D.B. and D.C. helped developing the methodology and interpret the results. P.M.V., N.W., M.R.R. and L.Y. performed sib-pairs analyses. K.E.K. and L.Y. performed quality control of UKB data. L.Y. and M.R.R. performed statistical analyses and simulations. L.Y. and P.M.V wrote the manuscript with the participation of all authors.

# Figure legends

2

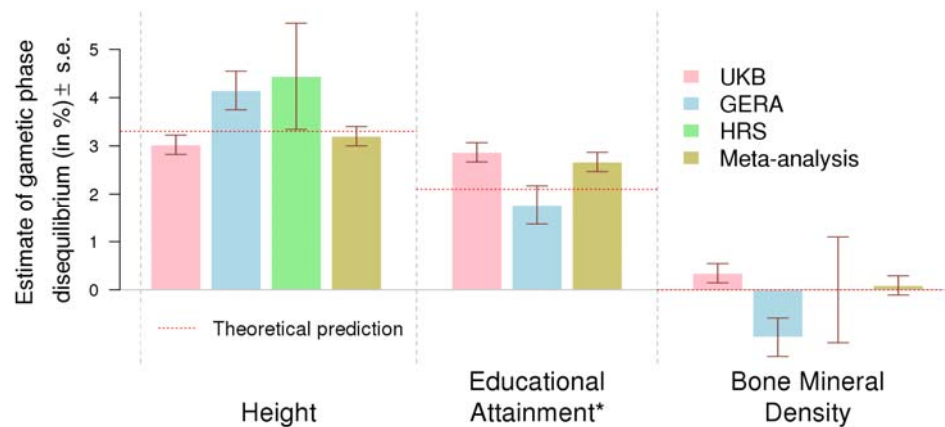


3

4

**Fig. 1.** Schematic illustration of the effect of assortative mating (AM) on the correlation between trait-associated alleles. Each line represents a chromosome of an individual in the population and each coloured bead represents an allele (orange: trait-increasing alleles (TIA); blue: trait-decreasing alleles) at a particular locus on that chromosome. Panel **a** represents a population undergoing AM and panel **b** represents a population undergoing RM. Under RM the distribution of alleles between odd and even chromosomes are uncorrelated (no-consistent pattern between chromosomes). Under AM, the distributions of alleles are correlated between chromosomes, such that the number of TIAs on odd chromosomes predicts the number of TIAs on even chromosomes.

15

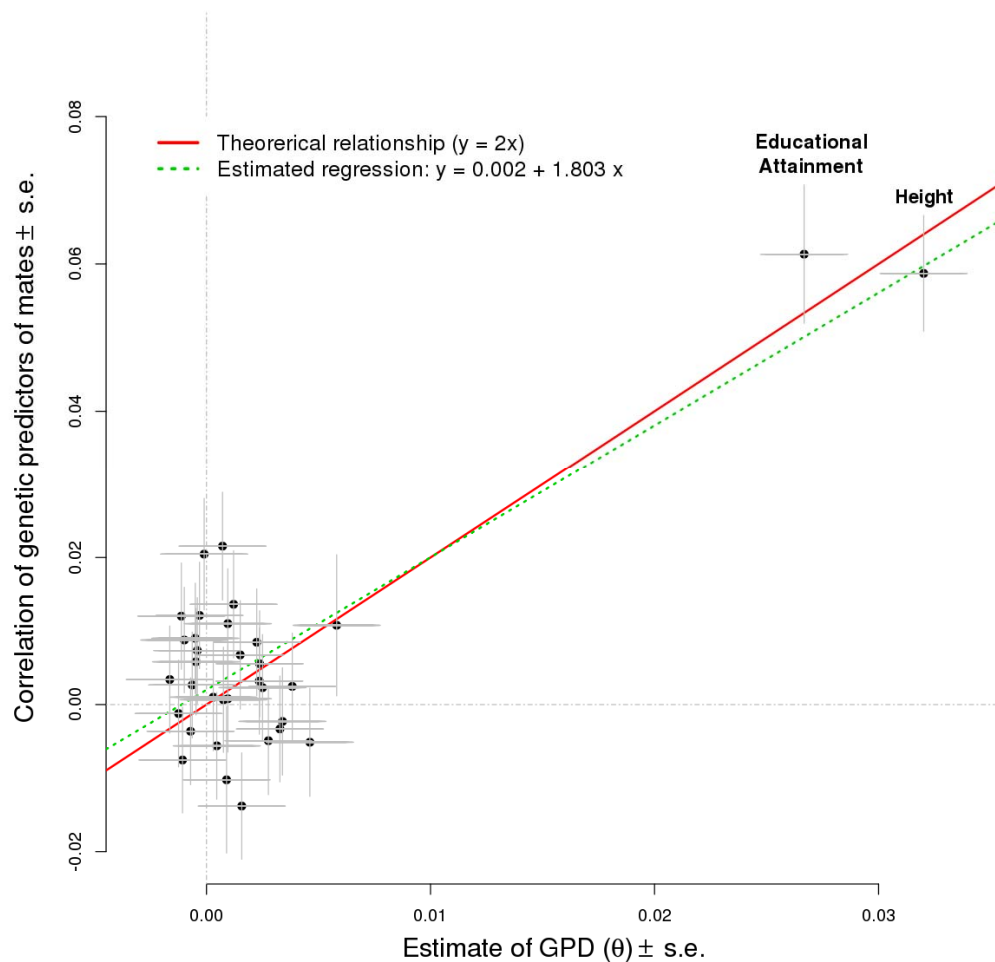


16

17 **Fig. 2.** Estimate of assortative mating (AM) induced gametic phase disequilibrium  
18 (GPD) among trait increasing alleles in three independent cohorts: UKB  
19 (N=348,502), GERA (N=53,991) and HRS (N=8,552). GPD is estimated as the  
20 correlation between trait-specific genetic predictors from SNPs on odd chromosomes  
21 versus even chromosomes. Bone mineral density was selected as a trait on which AM  
22 does not occur (negative control). Estimates are adjusted for 20 genotypic principal  
23 components from SNPs on either odd or even chromosomes to correct the effect of  
24 population stratification. \*The HRS cohort was not included in the meta-analysis of  
25 GPD estimates among educational attainment increasing alleles, as HRS was included  
26 in the Okbay *et al.* study. Theoretical predictions are obtained from equation (1).



27

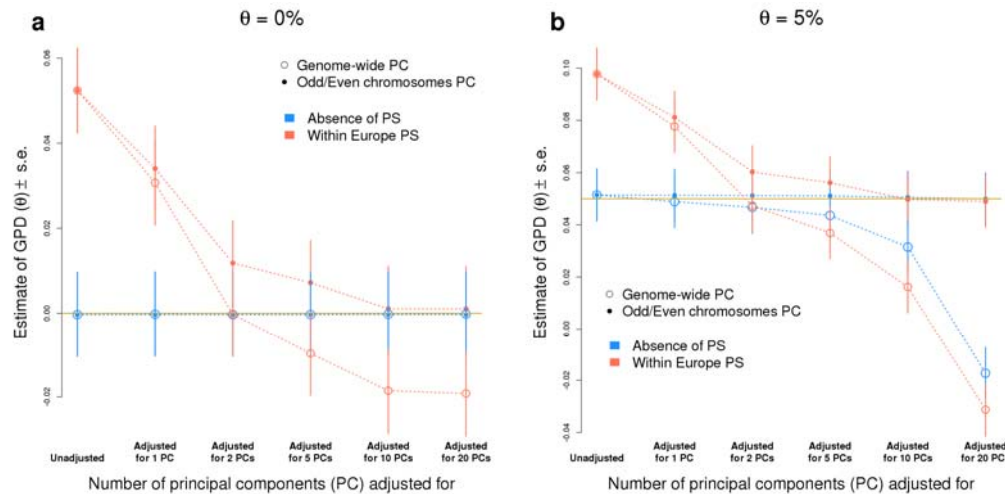


28

29 **Fig. 3.** Correlation of genetic predictors in 18,984 mates pairs (y-axis; values from  
30 **Table S2**) as a function of within-individual estimates of gametic phase  
31 disequilibrium (GPD: x-axis) for 32 complex traits and diseases (meta-analysis from  
32 **Table S1** in N=411,045 participants). Theory derived in **Supplementary Note 4**  
33 predicts a regression slope equal to 2. Estimated linear regression intercept and slope  
34 are 0.002 (standard error, s.e. 0.002) and 1.8 (s.e. 0.23) respectively.

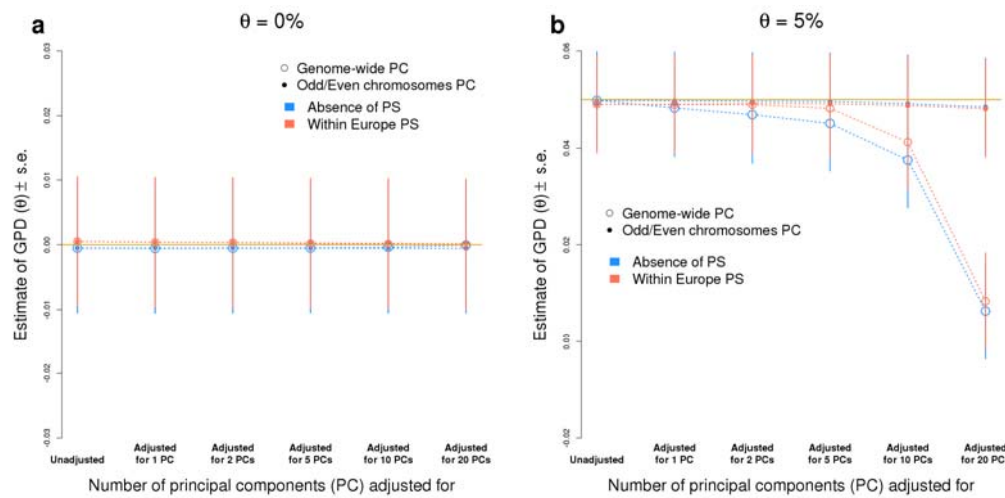
35

36



**Fig. S1.** Estimates of gametic phase disequilibrium (GPD), in simulated data (N=10,000) using allele frequencies of 697 height-associated SNPs, as a function of the number of genotypic principal components adjusted to correct for population stratification. Data were simulated assuming either no population stratification or within-Europe population stratification (**Supplementary Note 2**). In both cases, data were simulated under pure random mating ( $\theta=0$ , panel **a**) or under assortative mating ( $\theta=5\%$ , panel **b**). Estimates are obtained from unadjusted linear regression or adjusted for 1, 2, 5, 10 and 20 first principal components. Principal components were either calculated from SNPs on odd and even chromosomes (genome-wide principal components) or from SNPs on odd or even chromosomes separately. Standard Error (s.e.).

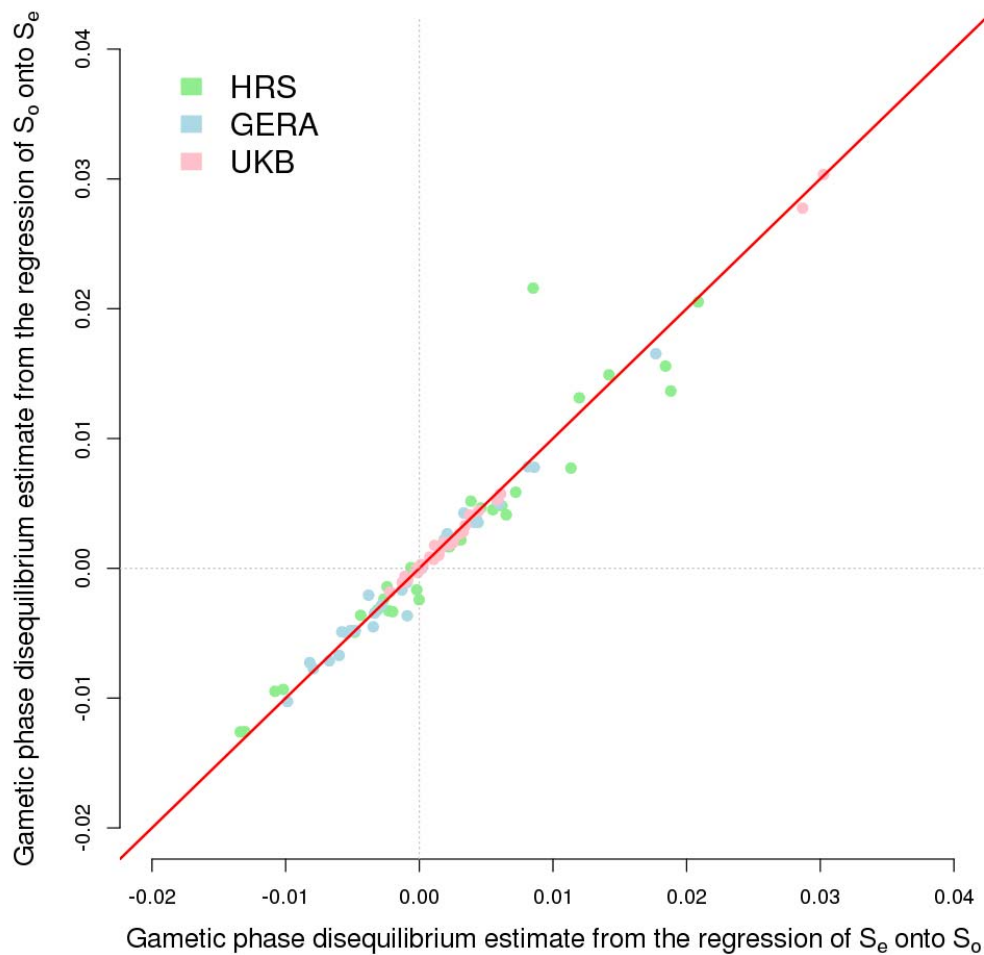
49



50

51 **Fig. S2.** Estimates of gametic phase disequilibrium (GPD), in simulated data  
52 (N=10,000) using allele frequencies of 1,000 randomly selected Hap Map 3 SNPs, as  
53 a function of the number of genotypic principal components adjusted to correct for  
54 population stratification. Data were simulated assuming either no population  
55 stratification or within-Europe population stratification (**Supplementary Note 2**). In  
56 both cases, data were simulated under pure random mating ( $\theta=0$ , panel **a**) or under  
57 assortative mating ( $\theta=5\%$ , panel **b**). Estimates are obtained from unadjusted linear  
58 regression or adjusted for 1, 2, 5, 10 and 20 first principal components. Principal  
59 components were either calculated from SNPs on odd and even chromosomes  
60 (genome-wide principal components) or from SNPs on odd or even chromosomes  
61 separately. Standard Error (s.e.).

62



**Fig. S3.** Comparison of estimates of gametic phase disequilibrium (GPD) from two approaches: (a) from the regression of genetic predictors from SNPs on even chromosomes ( $S_e$ ) onto of genetic predictors from SNPs on odd chromosomes ( $S_o$ ); and (b) from the regression of  $S_o$  onto  $S_e$ . The x-axis corresponds to approach (a) and y-axis to approach (b). The correlation between these two estimators is  $\sim 0.99$  across all three cohorts UKB ( $N=348,502$ ), GERA ( $N=53,991$ ) and HRS ( $N=8,552$ ).

Traits/Diseases	Number of SNPs to estimate $\theta$	$\theta_{\text{UKB}}$	$\theta_{\text{GERA}}$	$\theta_{\text{HRS}}$	$\theta_{\text{META}}$	Standard Error ( $\theta_{\text{META}}$ )	Meta-analysis $p$ -value	Heterogeneity $p$ -value
Height	9,447	3.02	4.15	4.44	3.20	0.16	$6.5 \times 10^{-89}$	0.03
Number of years of education	4,618	2.87	1.77	8.95*	2.67	0.19	$6 \times 10^{-46}$	0.02
Intelligence Quotient	3,356	0.61	0.33	1.42	0.58	0.19	$2.3 \times 10^{-3}$	0.65
Former smoker	2,011	0.44	0.82	-1.02	0.46	0.16	$3.8 \times 10^{-3}$	0.30
Attention Deficit Hyperactivity Disorder	2,902	0.58	-0.79	-0.31	0.38	0.16	0.02	0.01
Age at menopause	2,263	0.33	0.41	4.49*	0.34	0.16	0.03	0.85
Alcohol consumption (Continuous measurement)	2,186	0.19	0.86	0.31	0.28	0.17	0.09	0.39
Bipolar Disorder	2,358	0.38	-0.58	0.39	0.25	0.16	0.13	0.14
Age at menarche	3,385	0.18	0.38	1.56	0.24	0.16	0.13	0.42
Systolic Blood Pressure	1,687	0.22	0.44	-0.24	0.24	0.17	0.17	0.84
Crohn's Disease	2,669	0.14	0.60	1.20	0.22	0.16	0.17	0.45
Extraversion	1,666	0.30	-0.51	-1.30	0.16	0.16	0.33	0.10
Ever smoked	2,281	0.25	-0.60	0.72	0.15	0.16	0.35	0.18
Major Depressive Disorder	1,958	-0.11	-0.82	1.88	-0.16	0.18	0.35	0.09
Rheumatoid Arthritis	2,162	0.11	-0.09	0.85	0.10	0.11	0.38	0.47

Openness	1,674	-0.11	-0.54	2.09	-0.12	0.16	0.43	0.08
Type 2 Diabetes	2,904	0.18	-0.17	-0.33	0.12	0.16	0.46	0.72
Agreeableness	1,651	-0.01	-0.67	-0.44	-0.11	0.16	0.49	0.34
Anxiety case / control	1,872	-0.18	0.37	0.45	-0.10	0.14	0.49	0.38
Number of Cigarettes per Day	1,947	-0.13	0.02	-0.26	-0.11	0.17	0.52	0.95
Conscientiousness	1,688	0.16	-0.28	-0.06	0.10	0.16	0.55	0.63
Coronary Artery Disease	2,158	0.09	0.22	-1.26	0.08	0.16	0.63	0.44
Bone Mineral Density	5,859	0.34	-0.99	0.00	0.09	0.19	0.64	0.03
Autism Spectrum Disorder	1,811	0.11	-0.11	-0.17	0.07	0.16	0.65	0.88
Age at smoking onset	1,925	-0.03	-0.34	-0.23	-0.07	0.17	0.67	0.82
Anxiety (Factor score)	2,005	-0.09	0.01	0.46	-0.06	0.16	0.68	0.87
Post-Traumatic Stress Disorder	1,814	-0.09	0.19	0.22	-0.05	0.16	0.77	0.82
Body Mass Index (BMI)	2,362	0.01	-0.48	0.65	-0.04	0.15	0.78	0.43
Diastolic Blood Pressure	1,663	-0.01	0.58	-1.08	0.05	0.17	0.78	0.31
Waist/Hip ratio (adj. BMI)	1,940	0.01	0.27	-0.49	0.03	0.16	0.85	0.77
Alcohol consumption (Dichotomous measurement)	1,614	0.00	-0.38	1.13	-0.03	0.19	0.87	0.56
Schizophrenia	5,617	0.03	-0.34	0.62	-0.01	0.15	0.94	0.59

74

75 **Table S1.** Within-individual correlation ( $\theta$  in %) between trait-specific genetic predictors from SNPs on odd chromosomes versus even  
76 chromosomes for 32 traits and diseases. Estimates from UKB (N=348,502), GERA (N=53,991) and HRS (N=8,552) are denoted  $\theta_{\text{UKB}}$ ,  $\theta_{\text{GERA}}$  and  
77  $\theta_{\text{HRS}}$  respectively; and estimate from the meta-analysis of these three estimates is denoted  $\theta_{\text{META}}$ . \*As the HRS cohort was included in the  
78 educational attainment GWAS and the age at menopause GWAS,  $\theta_{\text{META}}$  calculations for these traits did not include data from HRS. All  
79 estimates are adjusted for 20 genotypic principal components from SNPs on either odd or even chromosomes. Standard errors (s.e.) of estimated  
80 correlations within each cohort are approximately inversely proportional to the square-root of the sample size (N). For UKB s.e. ~0.17%, for  
81 GERA s.e.~0.43% and for HRS s.e.~1.1%.

82

83

<b>Traits/Diseases</b>	$r_m$ (%)	s.e. ( $r_m$ )	$p$ -value	Number of mate pairs
<b>Height*</b>	<b>5.87</b>	<b>0.79</b>	<b><math>9.8 \times 10^{-14}</math></b>	<b>18,984</b>
<b>Number of years of education (EA)*</b>	<b>6.13</b>	<b>0.94</b>	<b><math>7.3 \times 10^{-11}</math></b>	<b>10,854</b>
Intelligence Quotient (IQ)	1.08	0.96	0.26	10,854
Former smoker	-0.51	0.74	0.48	18,984
Attention Deficit Hyperactivity Disorder	0.25	0.73	0.73	18,984
Age at menopause	-0.23	0.73	0.75	18,984
Alcohol consumption (Continuous measurement)	-0.49	0.73	0.5	18,984
Bipolar Disorder	0.23	0.72	0.75	18,984
Age at menarche	0.32	0.72	0.66	18,984
Systolic Blood Pressure	0.56	0.72	0.44	18,984
Crohn's Disease	0.85	0.73	0.24	18,984
Extraversion	-1.38	0.72	0.06	18,984
Ever smoked	0.68	0.74	0.36	18,984
Major Depressive Disorder	0.34	0.73	0.64	18,984
Rheumatoid Arthritis	1.10	0.75	0.14	18,984
Openness	-0.12	0.73	0.87	18,984
Type 2 Diabetes	1.36	0.73	0.06	18,984
Agreeableness	-0.75	0.72	0.29	18,984
Anxiety case / control	0.88	0.72	0.22	18,984
Number of Cigarettes per Day	1.20	0.72	0.96	18,984
Conscientiousness	0.08	0.72	0.91	18,984
Coronary Artery Disease	0.07	0.72	0.93	18,984
Bone Mineral Density (BMD)	-1.02	1.00	0.3	10,854
Autism Spectrum Disorder	2.15	0.73	0.003	18,984
Age at smoking onset	-0.36	0.73	0.62	18,984
Anxiety (Factor score)	0.27	0.72	0.71	18,984
Post-Traumatic Stress Disorder	0.58	0.72	0.42	18,984



Body Mass Index (BMI)	0.73	0.72	0.3	18,984
Diastolic Blood Pressure	-0.56	0.73	0.44	18,984
Waist/Hip ratio (adj. BMI)	0.10	0.73	0.89	18,984
Alcohol consumption (Dichotomous measurement)	1.21	0.72	0.09	18,984
Schizophrenia	2.05	0.76	0.006	18,984

84

85 **Table S2.** Between mates correlation (in %) of genetic predictors of 32 traits and  
86 diseases. Estimates are adjusted for 20 genotypic principal components from SNPs on  
87 both odd and even chromosomes. Significant correlations ( $p$ -value  $< 0.05/30$ ) are  
88 marked with a “\*”. Pairs involving UKB participants included in the EA, IQ and  
89 BMD GWAS were removed.

90

<b>Traits/Diseases</b>	<b>GWAS median sample size</b>	<b>LDSC heritability (h<sup>2</sup>)</b>	<b>LDSC heritability (standard error)</b>	<b>LDSC intercept</b>	<b>LDSC intercept (standard error)</b>	<b>LDSC Ratio statistic</b>	<b>LDSC ratio statistic (standard error)</b>
Number of years of education	293,723	12.7%	0.4%	0.93	0.01	-0.11	.
Intelligence Quotient	78,308	18.8%	1.0%	1.01	0.01	0.03	0.03
Height	252,083	34.2%	1.8%	1.23	0.03	0.12	0.02
Body Mass Index	233,692	13.5%	0.7%	0.66	0.01	-1.30	.
Waist/Hip ratio (adj. BMI)	142,438	9.0%	0.7%	0.86	0.01	-1.44	.
Attention Deficit Hyperactivity Disorder	55,374	23.2%	1.5%	1.03	0.01	0.10	0.03
Autism Spectrum Disorder	13,574	33.5%	4.4%	0.99	0.01	-0.17	.
Bipolar Disorder	16,731	47.1%	4.1%	1.01	0.01	0.03	0.06
Anxiety case / control	17,310	7.6%	3.0%	1.00	0.01	0.10	0.25
Anxiety (Factor score)	18,186	7.0%	2.7%	1.00	0.01	-0.05	.
Major Depressive Disorder	18,759	17.4%	2.6%	1.00	0.01	0.05	0.09

Post-Traumatic Stress Disorder	53,293	1.8%	0.8%	0.99	0.01	-0.48	.
Schizophrenia (SCZ)	150,064	24.3%	1.0%	1.04	0.01	0.05	0.02
Coronary Artery Disease	184,305	6.7%	0.5%	0.88	0.01	-0.88	.
Type 2 Diabetes (T2D)	152,599	7.8%	0.6%	1.00	0.01	-0.01	.
Crohn's Disease	20,883	50.7%	5.8%	1.02	0.01	0.09	0.05
Rheumatoid Arthritis	58,284	15.0%	2.6%	0.94	0.01	-0.50	.
Systolic Blood Pressure (SBP)	67,211	11.2%	0.9%	0.90	0.01	-1.77	.
Diastolic Blood Pressure (DBP)	67,201	11.0%	1.0%	0.91	0.01	-1.61	.
Age at menopause	70,000	13.4%	1.6%	0.99	0.01	-0.06	.
Age at menarche	182,416	11.9%	0.6%	0.97	0.01	-0.07	.
Conscientiousness	17,375	7.5%	3.2%	1.00	0.01	-0.06	.
Agreeableness	17,375	1.7%	3.0%	1.00	0.01	0.32	0.97
Openness	17,375	12.2%	2.8%	0.99	0.01	-0.40	.
Extraversion	17,375	4.5%	2.7%	1.00	0.01	-0.11	.
Ever smoked	74,053	7.5%	0.7%	1.00	0.01	-0.03	.
Age at smoking onset	74,053	2.0%	0.6%	1.00	0.01	-0.04	.
Number of Cigarettes per Day	74,053	3.1%	0.8%	1.00	0.01	0.09	0.12

Former smoker	74,053	3.5%	0.6%	1.00	0.01	0.01	0.13
Alcohol consumption (Binary)	70,460	0.8%	0.6%	0.99	0.01	-3.52	.
Alcohol consumption Continuous	70,460	5.2%	0.8%	1.01	0.01	0.17	0.09
Bone Mineral Density (BMD)	142,487	37.6%	4.1%	1.08	0.03	0.07	0.03

91

92 **Table S3.** Description of summary statistics from genome-wide associations (GWAS) on 32 traits analysed in this study. LD score regression  
93 (LDSC) was applied to each set of summary statistics to calculate SNP heritability and LDSC ratio statistics, a measure of population  
94 stratification. The LDSC software (version 1.0) does not calculate standard errors (s.e.) of the ratio statistic when it is < 0. In these cases, s.e. of  
95 the ratio statistics were replaced with “.”. URLs for downloading summary statistics used in this study are given in **Supplementary Note 1**.

96

## 97     **References**

- 98     1. Fisher, R. A. The correlation between relatives on the supposition of Mendelian  
99       inheritance. *Trans R Soc Edinb* 399–433 (1918).
- 100    2. Shine, R., O’connor, D., Lemaster, M. P. & Mason, R. T. Pick on someone your  
101       own size: ontogenetic shifts in mate choice by male garter snakes result in size-  
102       assortative mating. *Anim. Behav.* **61**, 1133–1141 (2001).
- 103    3. Jiang, Y., Bolnick, D. I. & Kirkpatrick, M. Assortative Mating in Animals. *Am.*  
104       *Nat.* **181**, E125–E138 (2013).
- 105    4. Pearson, K. & Lee, A. On the Laws of Inheritance in Man: I. Inheritance of  
106       Physical Characters. *Biometrika* **2**, 357–462 (1903).
- 107    5. Spuhler, J. N. Assortative mating with respect to physical characteristics. *Eugen.*  
108       *Q.* **15**, 128–140 (1968).
- 109    6. Mare, R. D. Five Decades of Educational Assortative Mating. *Am. Sociol. Rev.* **56**,  
110       15–32 (1991).
- 111    7. Silventoinen, K., Kaprio, J., Lahelma, E., Viken, R. J. & Rose, R. J. Assortative  
112       mating by body height and BMI: Finnish twins and their spouses. *Am. J. Hum.*  
113       *Biol. Off. J. Hum. Biol. Counc.* **15**, 620–627 (2003).
- 114    8. Stulp, G., Simons, M. J. P., Grasman, S. & Pollet, T. V. Assortative mating for  
115       human height: A meta-analysis: STULP et al. *Am. J. Hum. Biol.* (2016).  
116       doi:10.1002/ajhb.22917
- 117    9. Vandenburg, S. G. Assortative mating, or who marries whom? *Behav. Genet.* **2**,  
118       127–157 (1972).
- 119    10. Hippisley-Cox, J., Coupland, C., Pringle, M., Crown, N. & Hammersley, V.  
120       Married couples’ risk of same disease: cross sectional study. *BMJ* **325**, 636  
121       (2002).

- 122 11. Ajslev, T. A. *et al.* Assortative marriages by body mass index have increased  
123 simultaneously with the obesity epidemic. *Front. Genet.* **3**, 125 (2012).
- 124 12. Nordsletten, A. E. *et al.* Patterns of Nonrandom Mating Within and Across 11  
125 Major Psychiatric Disorders. *JAMA Psychiatry* **73**, 354 (2016).
- 126 13. Wright, S. Systems of mating. III. Assortative mating based on somatic  
127 resemblance. *Genetics* **6**, 144–161. (1921).
- 128 14. Crow, J. F. & Kimura, M. *An Introduction to Population Genetics Theory*.  
129 (Blackburn Press, 2009).
- 130 15. Nagylaki, T. Assortative mating for a quantitative character. *J. Math. Biol.* **16**,  
131 57–74 (1982).
- 132 16. Gimelfarb, A. Quantitative characters under assortative mating: gametic model.  
133 *Theor. Popul. Biol.* **25**, 312–330 (1984).
- 134 17. Bulmer, M. G. *The mathematical theory of quantitative genetics*. (Clarendon  
135 Press, 1980).
- 136 18. Sebro, R., Peloso, G. M., Dupuis, J. & Risch, N. J. Structured mating: Patterns  
137 and implications. *PLOS Genet.* **13**, e1006655 (2017).
- 138 19. Sebro, R., Hoffman, T. J., Lange, C., Rogus, J. J. & Risch, N. J. Testing for non-  
139 random mating: evidence for ancestry-related assortative mating in the  
140 Framingham heart study. *Genet. Epidemiol.* **34**, 674–679 (2010).
- 141 20. Risch, N. *et al.* Ancestry-related assortative mating in Latino populations.  
142 *Genome Biol.* **10**, R132 (2009).
- 143 21. Crow, J. F. & Felsenstein, J. The effect of assortative mating on the genetic  
144 composition of a population. *Soc. Biol.* **29**, 22–35 (1982).
- 145 22. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits:  
146 From Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).

- 147 23. Li, X., Redline, S., Zhang, X., Williams, S. & Zhu, X. Height associated variants  
148 demonstrate assortative mating in human populations. *Sci. Rep.* **7**, 15689 (2017).
- 149 24. Sampson, J., Kidd, K. K., Kidd, J. R. & Zhao, H. Selecting SNPs to Identify  
150 Ancestry. *Ann. Hum. Genet.* **75**, 539–553 (2011).
- 151 25. Wood, A. R. *et al.* Defining the role of common variation in the genomic and  
152 biological architecture of adult human height. *Nat. Genet.* **46**, 1173–1186 (2014).
- 153 26. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and  
154 Population-Based Linkage Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 155 27. Robinson, M. R. *et al.* Population genetic differentiation of height and body mass  
156 index across Europe. *Nat. Genet.* **47**, 1357–+ (2015).
- 157 28. Okbay, A. *et al.* Genome-wide association study identifies 74 loci associated with  
158 educational attainment. *Nature* **533**, 539–542 (2016).
- 159 29. Cesarini, D. & Visscher, P. M. Genetics and educational attainment. *Npj Sci.*  
160 *Learn.* **2**, (2017).
- 161 30. Bycroft, C. *et al.* Genome-wide genetic data on ~500,000 UK Biobank  
162 participants. *bioRxiv* 166298 (2017). doi:10.1101/166298
- 163 31. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: A Tool for  
164 Genome-wide Complex Trait Analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 165 32. Kemp, J. P. *et al.* Identification of 153 new loci associated with heel bone mineral  
166 density and functional involvement of GPC6 in osteoporosis. *Nat. Genet.* **49**,  
167 1468–1475 (2017).
- 168 33. Robinson, M. R. *et al.* Genetic evidence of assortative mating in humans. *Nat.*  
169 *Hum. Behav.* **1**, 0016 (2017).

- 170 34. Hugh-Jones, D., Verweij, K. J. H., St. Pourcain, B. & Abdellaoui, A. Assortative  
171 mating on educational attainment leads to genetic spousal resemblance for  
172 polygenic scores. *Intelligence* **59**, 103–108 (2016).
- 173 35. Conley, D. *et al.* Assortative mating and differential fertility by phenotype and  
174 genotype across the 20th century. *Proc. Natl. Acad. Sci.* **113**, 6647–6652 (2016).
- 175 36. Agrawal, A. *et al.* Assortative mating for cigarette smoking and for alcohol  
176 consumption in female Australian twins and their spouses. *Behav. Genet.* **36**,  
177 553–566 (2006).
- 178 37. Youyou, W., Stillwell, D., Schwartz, H. A. & Kosinski, M. Birds of a Feather Do  
179 Flock Together: Behavior-Based Personality-Assessment Method Reveals  
180 Personality Similarity Among Couples and Friends. *Psychol. Sci.* **28**, 276–284  
181 (2017).
- 182 38. Berg, J. J. & Coop, G. A population genetic signal of polygenic adaptation. *PLoS*  
183 *Genet.* **10**, e1004412 (2014).
- 184 39. Field, Y. *et al.* Detection of human adaptation during the past 2000 years. *Science*  
185 **354**, 760–764 (2016).
- 186 40. Tenesa, A., Rawlik, K., Navarro, P. & Canela-Xandri, O. Genetic determination  
187 of height-mediated mate choice. *Genome Biol.* **16**, 269 (2016).
- 188 41. Galinsky, K. J. *et al.* Fast Principal-Component Analysis Reveals Convergent  
189 Evolution of ADH1B in Europe and East Asia. *Am. J. Hum. Genet.* **98**, 456–472  
190 (2016).
- 191 42. Lachin, J. M. Introduction to sample size determination and power analysis for  
192 clinical trials. *Control. Clin. Trials* **2**, 93–113 (1981).
- 193 43. Allen, N. *et al.* UK Biobank: Current status and what it means for epidemiology.  
194 *Health Policy Technol.* **1**, 123–126 (2012).



- 195     44. 1000 Genomes Project Consortium *et al.* A global reference for human genetic  
196             variation. *Nature* **526**, 68–74 (2015).  
197