



King's Research Portal

Document Version
Peer reviewed version

[Link to publication record in King's Research Portal](#)

Citation for published version (APA):

Cheesman, R., Hunjan, A. K., Coleman, J. R. I., Ahmadzadeh, Y., Plomin, R. J., McAdams, T. A., Eley, T. C., & Breen, G. D. (2019). Comparison of adopted and non-adopted individuals reveals gene-environment interplay for education in the UK Biobank. Manuscript submitted for publication.

Citing this paper

Please note that where the full-text provided on King's Research Portal is the Author Accepted Manuscript or Post-Print version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version for pagination, volume/issue, and date of publication details. And where the final published version is provided on the Research Portal, if citing you are again advised to check the publisher's website for any subsequent corrections.

General rights

Copyright and moral rights for the publications made accessible in the Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognize and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Research Portal

Take down policy

If you believe that this document breaches copyright please contact librarypure@kcl.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.

Comparison of adopted and non-adopted individuals reveals gene-environment interplay for education in the UK Biobank

Rosa Cheesman¹, Avina Hunjan^{1,2}, Jonathan R. I. Coleman^{1,2}, Yasmin Ahmadzadeh¹,
Robert Plomin¹, Tom A. McAdams¹, Thalia C. Eley^{1,2*}, Gerome Breen^{1,2*}

¹Social Genetic & Developmental Psychiatry Centre, Institute of Psychiatry, Psychology &
Neuroscience, King's College London, UK

²NIHR Biomedical Research Centre for Mental Health; South London and Maudsley NHS
Trust, UK

*joint senior authors

Corresponding authors:

Prof. Thalia Eley (thalia.eley@kcl.ac.uk, +442078480863)

Prof. Gerome Breen (gerome.breen@kcl.ac.uk, +442078480409)

Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology &
Neuroscience, King's College London, London, SE5 8AF, UK

Running head: Adoption and gene-environment interplay for education

Word count (Introduction, Discussion, Acknowledgements): 1830

Abstract

Polygenic scores now explain ~10% of the variation in educational attainment. However, they capture not only genetic propensity but information about the family environment. This is due to passive gene-environment correlation, whereby the correlation between offspring and parent genotypes results in an association between offspring genotypes and the rearing environment. We measure passive gene-environment correlation using information on 6311 adoptees in the UK Biobank. Adoptees' genotypes are less correlated with their rearing environments, because they do not share genes with their adoptive parents. We find that polygenic scores are twice as predictive of years of education in non-adopted individuals compared to adoptees ($R^2 = 0.074$ vs 0.037 , $p = 8.23 \times 10^{-24}$). Individuals in the lowest decile of education polygenic score attain significantly more education if they are adopted, possibly due to educationally supportive adoptive environments. Overall, these results suggest that genetic influences on education are mediated via the home environment.

Keywords: educational attainment, polygenic scores, gene-environment interplay, adoption

Introduction

An important process by which genes and environments work together to influence behaviour is gene-environment correlation (Plomin et al. 2016). Gene-environment correlation refers to the association between the genotype an individual inherits from their parents and the environment in which they are raised (Plomin, DeFries, & Loehlin, 1977). Three forms of gene-environment correlation are typically distinguished: passive, active and evocative. An example of passive gene-environment correlation is that more educated parents are likely to provide both beneficial genes and educationally supportive family environments, such as books in the home, for their children. Therefore, shared genes confound associations between putative environmental variables and child attainment. Active and evocative gene-environment correlations reflect how genotypes lead to phenotypes: individuals select and evoke environments based on their genetically influenced traits.

It is essential to investigate gene-environment interplay in educational attainment, for several reasons. First, educational attainment is an important trait for individuals and society. Second, gene-environment correlation clearly matters for educational attainment. Adoption, twin, and instrumental variable research suggests that shared genes largely explain associations between parent and child attainment (Holmlund, Lindahl, & Plug, 2011). Third, polygenic scores, which index the genetic liability that each individual carries for a specific trait, are notably powerful for education attainment and now predict ~10% of the variation in years of education (Lee et al., 2018), with potential social implications (Plomin & von Stumm, 2018). However, it has only recently been shown that this prediction includes not only direct genetic effects on an individual's own education, but also indirect genetic effects

through relatives --i.e. predicting the family environment (Kong et al. 2018; Bates et al. 2018).

To disentangle causal processes affecting educational attainment, behaviour genetic study designs are needed. Adoption studies do this by removing overlapping genetic and environmental influences (passive gene-environment correlation). This is achieved by measuring the resemblance of adopted children with their birth parents, and with their adoptive parents. The former gives an estimate of direct genetic influence, independent of passively correlated environmental effects. The latter gives an estimate of shared environmental influence, free of correlated genetic effects. Passive gene-environment correlation may be estimated as the extent to which genes contribute more to the covariation between measures of the family environment and offspring traits in non-adoptive than adoptive families (Plomin, Loehlin, & DeFries, 1985). Notably, other forms of gene-environment correlation are still present in adoptees, since heritable proclivities lead them to select and evoke experiences.

More recently, researchers have applied genomic tools to family data to estimate direct and indirect effects on educational attainment (Bates et al., 2018; Domingue, Belsky, Conley, Harris, & Boardman, 2015; Kong et al., 2018; Selzam et al., 2019; Wertz et al., 2018; Young et al., 2018). These designs are conceptually related to adoption designs, since they account for shared genes between parents and offspring. For example, genetic variants that were not passed on by parents can only have indirect effects on offspring traits, through genetically-influenced parental behaviour (Bates et al., 2018; Kong et al., 2018). When controlling for indirect effects with an education polygenic score based on non-transmitted variants, the variance explained by the transmitted score shrank from 5 to 2% (Kong et al., 2018). The

non-transmitted score also independently predicted attainment. The family environment is an important contributor to polygenic score prediction because it is adding to estimates of genetic influence, and because parents still influence their offspring after controlling for shared (transmitted) genes.

This study draws on an unusually large, and relatively unexplored, sample of adoptees, and harmonises a traditional quantitative genetic approach with modern genomic tools. Our main aim was to use the ‘natural experiment’ created by adoptive placement to measure the importance of passive gene-environment correlation for educational attainment. When children are adopted by non-relatives, the indirect genetic path between the rearing environment and their traits is not present because adoptive parents are not genetically related to adopted children. Three hypotheses follow. First, the phenotypic variance should be lower in adoptees compared to non-adopted individuals, because adoptees do not have the additional source of variance of passive gene-environment correlation (Loehlin & De Fries, 1987; Plomin, 1994). It could also be because adoptive families may vary less in socio-economic status or may be selected for perceived parenting ability (Natsuaki et al., 2019; Rutter, 2006). Second, if passive gene-environment correlation inflates heritability estimates, then heritability should be lower in adoptees than in non-adopted individuals, because adoptees are reared in environments that are less correlated with their genotypes. Third, for the same reason, the variance explained by polygenic scores will be lower in adoptees, and may be closer to the direct genetic effect of an individual’s own DNA.

Method

Sample, genotype quality control and phenotype definition

The UK Biobank is an epidemiological resource including British individuals aged 40 to 70 at recruitment (Allen, Sudlow, Peakman, Collins, & UK Biobank, 2014). UK Biobank participants were asked “Were you adopted as a child?”. 8,040 individuals said yes, and 541,889 individuals said no. No additional information was collected on factors that are understood to reduce the representativeness of adoptees as a study sample: the age of adoption, whether the adoption was domestic or international, or whether individuals were adopted by biological relatives. Genome-wide genetic data came from the full release of the UK Biobank data, and were collected and processed according to the quality control pipeline (Bycroft et al., 2018). We restricted analyses to individuals with full phenotypic data for education, who also passed genotype quality control criteria. This left 6,311 adopted and 375,343 non-adopted individuals for analysis.

Genotype quality control criteria were: common genetic variants of minor allele frequency > 0.01 that were directly genotyped or imputed with high confidence (INFO metric > 0.4); and individuals with genotype call rate $> 98\%$ who had concordant phenotypic and genetic gender information and who were unrelated to others in the dataset (less than third degree relatives). We performed removal of relatives using a “greedy” algorithm to minimise the exclusion of adoptees. To reduce confounding from population stratification, all analyses were limited to individuals of European ancestries, as defined by 4-means clustering on the first two genetic principal components provided by the UK Biobank. We also controlled for 10 ancestry principal components of the European sample in all genomic analyses.

Years of education, a proxy for educational attainment, was defined according to ISCED categories, as in previous genomic studies of the phenotype in UK Biobank and other samples (Lee et al., 2018). The response categories were: none of the above (no qualifications) = 7 years of education; Certificate of Secondary Education (CSEs) or equivalent = 10 years; O levels/GCSEs or equivalent = 10 years; A levels/AS levels or equivalent = 13 years; other professional qualification = 15 years; National Vocational Qualification (NVQ) or Higher National Diploma (HNC) or equivalent = 19 years; college or university degree = 20 years of education.

Statistical analyses

Phenotypic comparisons. First, we formally tested the hypothesis that non-adopted individuals show greater phenotypic variance than adopted individuals due to the presence of an additional source of variance (passive gene-environment correlation). A non-parametric test was used given the non-normal distribution of the education years variable (Brown & Forsythe, 1974). This test is based on absolute deviations from the median, rather than the group mean. We also tested for differences in education years, age, and sex between the two groups, using a Wald test, z-test, and Wilcoxon test, respectively.

SNP heritability estimation. Second, to test the hypothesis that heritability is lower in adoptees, whose rearing environments are less correlated with their genotypes, we estimated the variance explained by common genetic variants for years of education in adoptees using Genomic-RElatedness-based restricted Maximum-Likelihood (GREML) (Yang, Lee, Goddard, & Visscher, 2011), and compared this to the heritability estimate for non-adopted individuals. The method estimates heritability as the extent to which genetic similarity among unrelated individuals can predict their trait similarity. In GREML, a matrix

of genomic similarity for each pair of unrelated individuals across genotyped variants is compared to a matrix of their pairwise phenotypic similarity using a random-effects mixed linear model, such that the variance of a trait can be decomposed into genetic and residual components, using maximum likelihood. We used two genetic relatedness matrices: one for adopted individuals, and a second for a subset of 6,500 non-adopted individuals. This was to enable comparison of two similarly sized samples, and to reduce the computational burden that results from scaling GREML to a sample as large as the UK Biobank. For both genomic matrices we used a relatedness cutoff of 0.025. Sub-samples were made using the “sample_n” function in the dplyr package in R (version 3.5). We compared these results to heritability estimates derived from a second method, LD score regression (LDSC) (B. K. Bulik-Sullivan et al., 2015). Unlike GREML, LDSC does not require individual-level data, allowing it to be computationally feasible to estimate the heritability of education in the full sample of non-adopted individuals. LDSC also enabled us to estimate genetic correlations (see below).

Polygenic scoring. Third, we tested whether the power of polygenic scores is greater for individuals who were reared with their biological relatives than for adoptees. The sample of non-adopted individuals was subdivided into three independent groups for polygenic score analyses. Our first sample consisted of 318,843 non-adopted individuals for genome-wide association analysis (GWA). The purpose of this was to estimate the effect sizes of associations between genome-wide genetic variants and years of education, to use for the creation of individual-level polygenic scores. We derived our base summary statistics file for years of education by meta-analysing summary statistics from our own GWA analysis in this subsample with independent summary statistics obtained from the Social Science Genomics Consortium (excluding UK Biobank and 23&Me) (Lee et al., 2018). The sample size for

these external summary statistics was 324,160, leading to a total sample size of individuals in our GWA meta-analysis of 643,003.

Our second independent sample included 50,000 individuals to use for training our polygenic scores for years of education, i.e. identifying the optimal p-value threshold for inclusion of SNPs. The standard set of P-values in PRSice 2 were tested: 0.001, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 1 (Choi Shing Wan, n.d.).

Our third independent sample included 6,500 individuals, to match the sample size of adopted individuals, in which to run polygenic prediction models. In these prediction models we regressed the years of education phenotype in the UK Biobank on polygenic scores for years of education in adoptees, and then repeated the analysis in the 6,500 non-adopted individuals. In this third set of analyses we used a set p-value threshold obtained from the training step. This exact sample was the same as the one used to estimate SNP heritability of years of education. Notably, this polygenic score analysis is better-powered than the SNP heritability analysis, since it capitalises on the power of the large discovery sample (N=643,003).

As a negative control, we tested the polygenic prediction comparison between adoptees and non-adopted individuals for height, which has not shown evidence of passive gene-environment correlation in previous studies (Kong et al. 2018; Selzam et al. 2019). As with the education analysis, we trained the polygenic score based on the largest independent association study (Wood et al., 2014) in the sample of 50,000 individuals, and then tested the prediction at the best p-value threshold in our two independent and similarly sized samples of adopted and non-adopted individuals.

Supplemental analyses.

Heritability of adoption status. Substantial heritability of our environmental moderator might affect the interpretation of our main results. To explore this, we also estimated the heritability of adoption status using LD score regression in the full sample (N=381,654 individuals). The genetic 'influence' on adoption largely arises in the biological parent generation because heritable traits influence the likelihood of adoption of their child.

Polygenic score by adoption interaction analyses. Differences in genetic influences on the same trait across contexts - in this case adoption - can also be conceptualised as gene-environment interaction, whereby the impact of genes on educational attainment may be contingent on adoption status. We aimed to further explore our main results by testing a formal polygenic score by adoption interaction regression model. The model included main effects for polygenic score for years of education, adoption, and covariates, plus the interaction term as well as interaction terms for polygenic score and adoption with each covariate (Keller, 2014). We tested a linear model for additive interaction and a logistic model for multiplicative interaction. To visualise any interaction, we plotted the regression slopes for polygenic prediction of educational attainment for adopted and non-adopted individuals (with both variables standardised to have a mean of 0 and a standard deviation of 1). Additionally, we stratified polygenic scores for adopted and non-adopted individuals overall (N=12,811) into deciles and tested for mean differences in years of education between adopted and non-adopted groups in each decile.

Qualitative differences in the genetic influence on education by adoption status. We assessed whether education is driven by the same set of genetic influences in adopted and

non-adopted individuals. First, we estimated the genetic correlation between education in our samples. For this, we ran genome-wide association analyses of years of education in the full sample of non-adopted individuals (N= 375,343) and in the sample of adoptees, then estimated the genetic correlation between them using LD score regression. Second, we tested whether education is genetically linked to different traits between adoptees and non-adopted individuals. To this end, we estimated genetic correlations between education years and 247 traits available on LD Hub, for both adopted and non-adopted individuals. We compared the magnitudes of genetic correlations between education years and other variables between the adopted and non-adopted samples with z-tests.

Birth year-related differences in genetic influence. During the period when UK Biobank participants were growing up (1930s-70s), access to education increased, and there was great change in the norms and regulations surrounding reproduction, contraception and adoption. Previous studies have found that genetic influence on years of education increased in this period in the UK, since environmental differences between people had less influence on whether they stayed in education (Lee et al., 2018). We investigated temporal change in patterns of genetic influence on education in adopted versus non-adopted individuals by stratifying polygenic prediction analyses according to year of birth. Specifically, all individuals were split into 7 mutually exclusive birth-year groups, each with a range of 5 years, and polygenic score analyses were conducted separately for each of the year groups.

All analyses (SNP heritability, polygenic scoring, GWA) controlled for the following covariates: sex, age, 10 ancestry principal components, and factors capturing genotyping batch and centre. The majority of the analyses were completed in R version 3.5. GREML was performed in the GCTA software (Yang et al., 2011). Genome-wide association meta-

analysis was performed in METAL (Willer, Li, & Abecasis, 2010). Polygenic score analyses were performed in PRSice 2 (Choi Shing Wan, n.d.). To compare polygenic score results between adopted and non-adopted individuals, we obtained bootstrapped standard errors for the R^2 statistics using the boot package in R, with 1000 replications. Genome-wide genetic correlations were estimated using LDSC (B. Bulik-Sullivan et al., 2015) and LD Hub (Zheng et al., 2017). The UK Biobank is a controlled-access public dataset available to all bona fide researchers.

Results

Sample analysed

The total sample of individuals with education phenotype data and quality-controlled genotype data was 381,654. As described in the Methods, individuals were split into 4 mutually-exclusive groups: a) adopted as children (N=6,311), b) 318,843 non-adopted individuals for genome-wide association analysis, c) 50,000 non-adopted individuals for training of polygenic scores, and d) 6,500 non-adopted individuals for genomic analyses comparing to adoptees. Non-adopted individuals were randomly placed into groups b, c and d.

Phenotypic results

Phenotypic differences between adoptees and non-adopted individuals were generally modest in size but, due to the large sample size in this study, several were statistically significant (see Table 1). We found that non-adopted individuals showed significantly greater variance in their years of education than adoptees (26.2 vs 25.8; $p=0.002$ compared to 6500 non-adopted individuals in group d; $p=3.2 \times 10^{-5}$ compared to all non-adopted individuals). Table 1 gives descriptive statistics for education years, age and sex in the two groups. Adoptees in the UK Biobank were significantly younger on average ($p=0.026$ compared to group d; $p=0.009$ compared to all non-adopted individuals) although point estimates were similar (56.4 versus 56.7). There were significantly more males in the adopted group ($p=0.033$ compared to group d; $p=0.008$ compared to all non-adopted individuals), but the magnitude of the difference is small (48 versus 46% male). Adoptees had significantly fewer years of education ($p=3.3 \times 10^{-11}$ compared to group d; $p < 2.2 \times 10^{-16}$ compared to all non-adopted individuals in the UK Biobank). This is also reflected in the lower percentage of college attendees (20 years of

education in Table 1) in the adopted group (28% compared to 33%). All comparison results were consistent between the large and small samples of non-adopted individuals.

Table 1: Comparative analysis of phenotypes in adoptees versus non-adopted individuals.

		Adopted (N=6311)	Non-adopted (N=375343)	Non-adopted (group d; N=6500)
	Age	56.4 (8.53)	56.7 (8.01)	56.7 (8.06)
	Sex	48% male (N=3010)	46% male (N=172706)	46% male (N=2978)
Education				
Years	7	1209 (19%)	62651 (17%)	1064 (16%)
	10	1780 (28%)	100210 (27%)	1709 (26%)
	13	749 (12%)	43448 (12%)	755 (12%)
	15	350 (6%)	19428 (5%)	354 (5%)
	19	433 (7%)	24300 (6%)	428 (7%)
	20	1790 (28%)	125306 (33%)	2190 (34%)

Note: Adoptees were compared to the full sample of non-adopted individuals, and to our smaller sub-sample used for genomic analyses (group d).

SNP heritability estimates

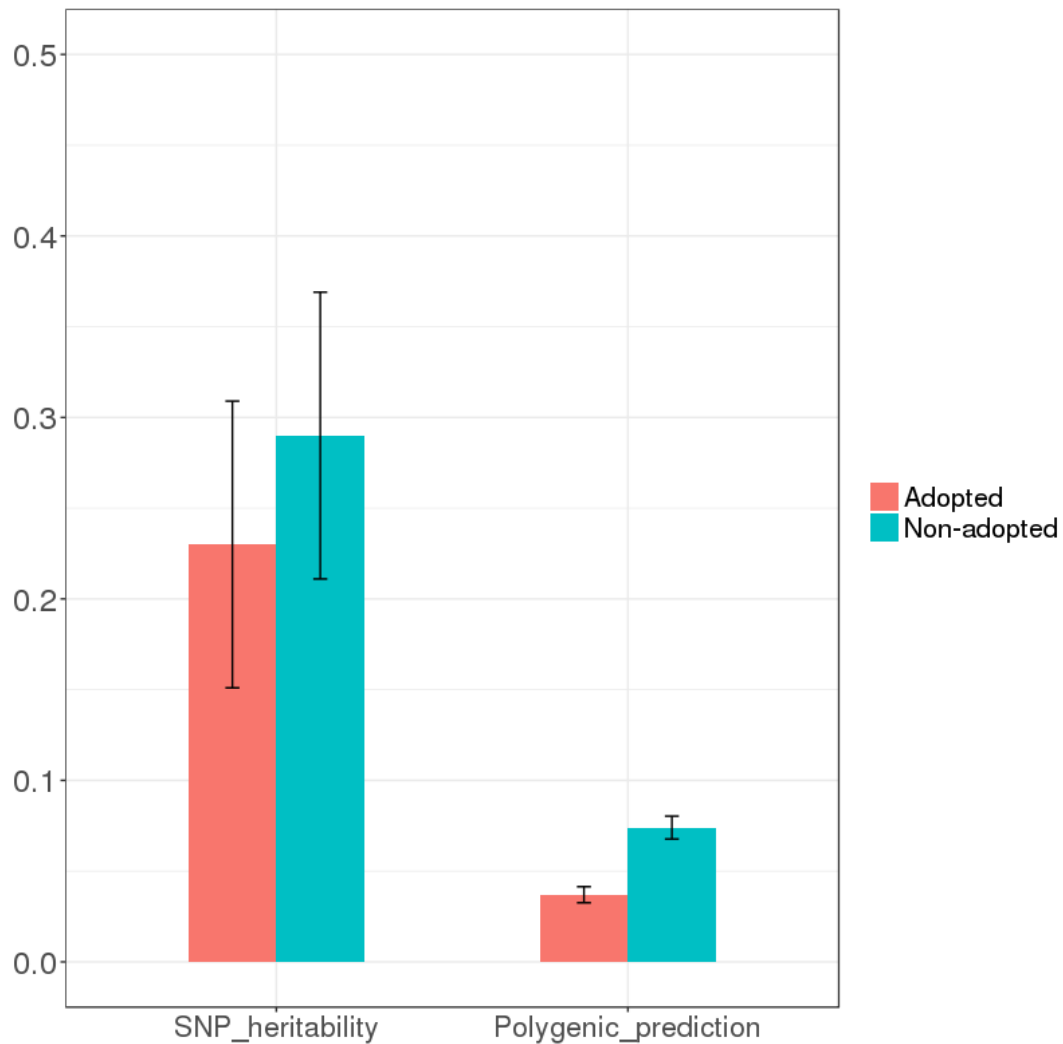
Figure 1 compares GREML-derived SNP heritability estimates for years of education in adopted individuals versus non-adopted individuals (left-hand bars). The estimate of heritability was larger in individuals reared with their relatives (0.29 [se = 0.079]) compared to adopted individuals (0.23 [se = 0.079]). However, confidence intervals were wide and overlapped, so the difference in heritability was not significant.

It was not computationally feasible to estimate the heritability of education using all non-adopted individuals with GREML. Notably, though, the LD score regression-derived heritability was 0.17 (se = 0.005) in the full sample of non-adopted individuals (N= 375,343), and 0.14 (se=0.073) for adoptees, corroborating the pattern of results found using GREML. LD score regression estimates are typically lower than GCTA-GREML-derived estimates (Evans & Keller, 2018).

Polygenic prediction results

Figure 1 also shows that twice as much phenotypic variance in years of education was explained by polygenic scores for education years in non-adopted individuals (0.074) as in adoptees (0.037). This difference was highly significant ($p= 8.23 \times 10^{-24}$). The optimal threshold for inclusion of SNPs was $p=1$ (Table S1). Table S2 shows the full results from the polygenic prediction analyses.

Figure 1: Estimates of the variance explained by common SNPs for years of education, and of the variance explained by polygenic scores for education polygenic scores, in adoptees compared to individuals who were reared with their relatives, plus 95% confidence intervals.



Note: sample sizes for polygenic prediction analyses were 6,311 and 6,500 for adopted and non-adopted individuals respectively; sample sizes for GREML heritability analyses were lower (6,227 for adopted and 6,362 for non-adopted individuals) since relatives were removed at a cutoff of >0.025 . For polygenic score results, CIs were obtained by bootstrapping with 1000 replications.

For our negative control analysis of height, as expected, the variance explained by the polygenic score in adoptees (0.127, se = 0.008) versus non-adopted individuals (0.134, se = 0.008) was not significantly different ($p=0.62$). The optimal threshold for inclusion of SNPs in the polygenic score was $p=0.001$.

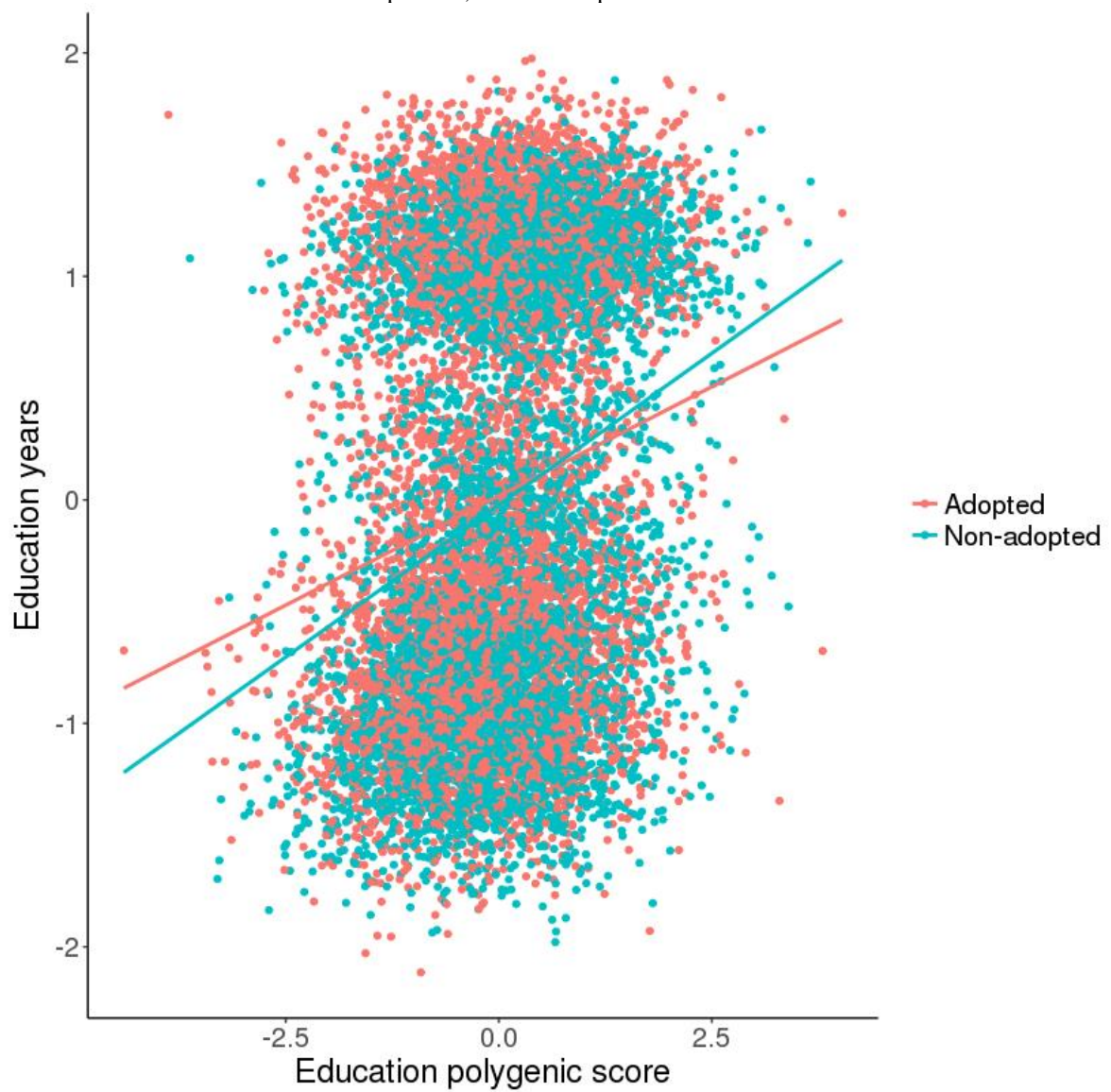
Supplemental analyses

Heritability of being adopted. We found a liability scale SNP heritability of being adopted of 0.059 (se = 0.004), assuming the population prevalence of adoption is identical to the sample prevalence (1.7%). If the actual population prevalence differed and was, for example, 0.7% or 2.7%, the liability scale SNP heritability would become 0.047 (se = 0.002) or 0.066 (se = 0.005), respectively. Adoption status showed significant genetic correlations with education, age at first birth, depression and obesity after correcting for multiple testing (see Table S5), although these correlations should be viewed with caution given the low SNP heritability of adoption. Adoption status could be significantly predicted by the education years polygenic score ($R^2=0.008$, $p < 2 \times 10^{-16}$). The heritability of adoption is low but may confound our between-group comparisons.

Polygenic score by adoption interaction. We tested a formal interaction model to further examine the finding that genetic influences on education are weaker in the sample of adoptees. The interaction between polygenic score and adoption status in predicting years of education is visualised in Figure 2. The regression slope is significantly steeper in the non-adopted group, indicating that years of education increases more as education polygenic scores increase in this group compared to the adopted group. See Table S3 for the full interaction model results. Using linear regression, we confirmed that polygenic prediction of education interacts beyond additivity with adoption status (interaction estimate = -0.33; $p=2.66 \times 10^{-4}$). This means that polygenic scores had a smaller association with education in adoptees. Then using logistic regression instead of linear regression, we also found that the interaction exceeded multiplicativity (interaction estimate = -0.18; $p = 0.0009$). The finding of interaction exceeding both additive and multiplicative models means that the combined effect of education polygenic score and adoption status is not scale dependent and is greater than either the sum or product of their individual effects, respectively.

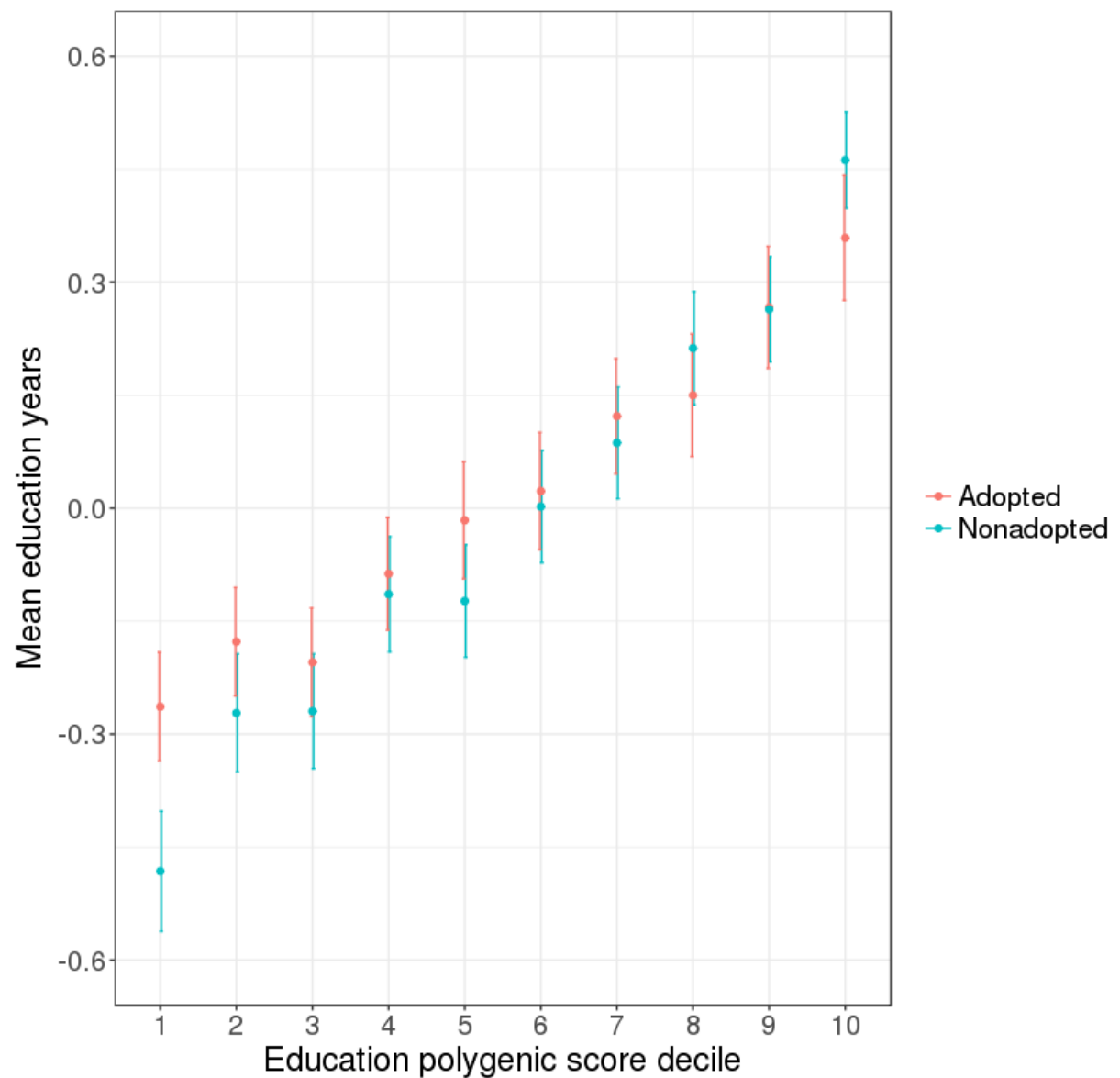
To further explore the interaction, we plotted the mean years of education per decile of polygenic score for education years, for adopted and non-adopted individuals. Figure 3 shows that for individuals in the lowest decile of education polygenic score, those who were adopted as a child achieved a substantially higher mean years of education (standardised) compared to non-adopted individuals (-0.24 [se = 0.03] versus -0.40 [se = 0.03]). This mean difference between adopted and non-adopted individuals was significant in the bottom decile ($p=7.05 \times 10^{-5}$), but not for other deciles of polygenic load. See Table S4 for full results of the decile analysis.

Figure 2: Regression of years of education on polygenic score for education, comparing 6311 adoptees to a sample of 6,500 non-adopted individuals.



Note: the two clusters of data-points reflect the distinct groups of individuals who did and did not attend university.

Figure 3: Mean education years (standardised) per decile of polygenic score for education years (also standardised), for adopted and non-adopted individuals, plus 95% confidence intervals.



Qualitative differences in genetic influences according to adoptee status. We

found that largely the same genetic influences are operating on education regardless of adoption status. First, the genetic correlation between education years in adopted and non-adopted individuals was not significantly different from 1 (0.81 [se = 0.21]). Second, we found no evidence that educational attainment is associated with different traits in individuals

who were adopted. Figure S1 presents estimates of genetic correlations between years of education and 247 external traits, comparing the adopted and non-adopted samples. None of these were significantly different between adoptees and non-adopted individuals after multiple testing correction. Due to the relatively small sample of adopted individuals, these results should be interpreted with caution.

Year-of-birth stratification analysis. Our final sensitivity analysis tested for differences according to year of birth in polygenic prediction from direct effects (indicated by the variance explained in the adoptees) versus from passive gene-environment correlation (indicated by the difference in variance explained between non-adopted individuals and adoptees). We found small, non-significant differences in the variance explained by polygenic scores for education depending on the year-of-birth group considered. Figure S2 shows that polygenic prediction remains generally stable for the adoptees across generations at ~ 0.04 , and any differences between age strata were non-significant. We note that sub-sampling reduced the statistical power to detect differences within and between groups across time. See Table S6 for sample sizes of each year-of-birth group.

Discussion

Cumulatively, our findings suggest that the family environment provided by relatives plays an important role in the manifestation of genetic effects on education. The educational attainment of individuals who were adopted away from their parents as children had significantly less variance explained by polygenic scores ($R^2 = 0.04$ versus 0.07 ; difference test $p = 8.23 \times 10^{-24}$). The variance explained by polygenic scores in years of education in adoptees (0.04) approximates the prediction coming from the direct effects of individuals' own DNA. The difference between the variance explained in non-adopted individuals and adoptees suggests that about half of the predictive power of polygenic scores for educational attainment comes from passive gene-environment correlation. We also found that individuals in the lowest decile of polygenic score attained significantly more years of education if they were adopted.

By showing that polygenic scores for education are twice as powerful in non-adopted individuals compared to adoptees, we suggest that genetic influence on educational attainment is magnified when individuals are reared by their close genetic relatives, with whom they share both genes and environments. Our results agree with recent evidence showing that the effects of passive gene-environment correlation reduced the variance explained by polygenic scores by 30-50% (Kong et al., 2018; Selzam et al., 2019).

Notably, in line with other recent research (Kong et al., 2018; Selzam et al., 2019), we use the term 'direct genetic effect' to refer to the effect of the polygenic score among adoptees, controlling for passive gene-environment correlation. However, it cannot be assumed that estimating the direct effect of a polygenic score is tantamount to isolating a straightforward

“purely genetic” effect or a “genetic propensity”. Genetic effects are never truly direct, but are always behaviourally mediated and expressed in the context of an environment. Active and evocative gene-environment correlation mechanisms are essential in how genes influence traits in everyone, including adoptees (Plomin, 2014), and these are included in estimates of direct genetic influence.

Our observation that individuals in the lowest decile of education polygenic score attain significantly more education if they are adopted could be due to educationally supportive adoptive environments. This agrees with previous evidence showing that adoptees had higher school achievement and intelligence test scores than non-adopted siblings or peers who stayed with their birth family (van Ijzendoorn, Juffer, & Poelhuis, 2005), and that such advantages are retained in their adult qualifications (Maughan, Collishaw, & Pickles, 1998). The specificity of this results to adoptees in the lowest polygenic score decile links to previous evidence that the ‘boosting’ effect may be stronger in higher socio-economic status adoptive families, and for children rescued from poverty (Duyme et al. 1999; Turkheimer, 1991). This environmental effect of adoptive parents might suggest that efforts to help individuals stay in education can be effective for those with less genetic propensity for education.

These results should be viewed in light of several limitations. First, interpreting genetic influence in adoptees as direct and free of passive gene-environment correlation requires that close biological relatives were not involved in the education of the adoptees. Unfortunately, the UK Biobank contains no information about the age of adoption beyond that it occurred in childhood, nor whether individuals were adopted by relatives or were able to identify and contact their biological parents. This knowledge would have allowed us to exclude

individuals who were not solely socialised with adoptive families, and therefore to make a precise comparison to individuals who were reared with their birth parents. However, polygenic prediction of education still differed markedly between the two groups, even though adoptees may have been in contact with their biological relatives. Thus, the effects of passive gene-environment correlation may contribute even more than half of the predictive power of education polygenic scores, as we estimated here.

A second caveat is lack of generalisability. The UK Biobank is not representative of the general population, since there is ‘healthy and wealthy’ volunteer selection (Fry et al., 2017; Keyes & Westreich, 2019), and we have only analysed data on individuals with European ancestries. Adoptees may not be random samples of the population. Indeed, adoption status is not a purely random environmental exposure but is slightly heritable (our estimate is 6%), and this may confound our between-group comparisons. Moreover, due to the lack of detailed adoption data in the UK Biobank, we were unable to conduct sensitivity analyses adjusting for aspects of adoption that have known associations with important outcomes including school performance. For example, it would have been useful to know whether adoptions were domestic or international, or in the context of childhood adversity and/or institutionalisation (Howard et al. 2004).

Similarly, adoptive parents tend to differ systematically from other parents: they are likely to be more educated, more socially advantaged, and to have lower rates of psychopathology (Rutter, 2006). A recent US study found that children are adopted into households that differ in average parental education compared to biological children (Domingue & Fletcher, 2019). This probably applies to the UK Biobank, although we cannot be certain, due to the lack of parental data. If adoptive families are more homogeneous with respect to these

characteristics, environmental variance may contribute less to differences in educational attainment among adoptees, and trait heritability estimates are consequently likely to be higher. However, the fact that lower environmental variance may act to *inflate* genetic influence in adoptees compared to non-adopted individuals makes our finding of significantly higher polygenic prediction in non-adopted individuals more striking. Again, the effects of passive gene-environment correlation for education may be even greater than we estimate.

There are several advantages of using the present adoption design for distinguishing direct genetic influence from passive gene-environment correlation. Unlike other methods, our approach does not require intergenerational data, which is valuable but has its own issues, such as cohort differences in genetic effects. Analysing the adoptees in the UK Biobank also bypasses several limitations of traditional adoption studies, including low sample size, and reliance on weak indirect proxies for inherited load for specific traits (birth parent trait status rather than individual-level polygenic scores). However, future progress in understanding the mechanisms driving the transmission of educational attainment will require intergenerational, longitudinal, genetically informative datasets, including detailed characterisation of the home environment. A developmental approach is useful, since gene-environment correlations likely arise early in childhood, and there will be complex reciprocal effects across time. Researchers have already started to pinpoint genetically-influenced aspects of families that are associated with children's education polygenic scores (Krapohl et al., 2017; Wertz et al., 2018, 2019).

The evidence presented in this study highlights the importance of the family environment to causal mechanisms influencing individual differences in educational attainment. These can be through possessing genes that shape the educational environment provided for offspring that

also directly influence attainment in the child, or through providing an educationally supportive environment for your adopted child.

Acknowledgments

We thank the scientists involved in the construction of the UK Biobank and all of the participants who have shared their life experiences with investigators in the UK Biobank. This research has been conducted using the UK Biobank Resource, under the application 18177 (with thanks to Paul F. O'Reilly). This study represents independent research part funded by the National Institute for Health Research (NIHR) Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and King's College London. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care. High performance computing facilities were funded with capital equipment grants from the GSTT Charity (TR130505) and Maudsley Charity (980). T.C. Eley is part funded by a program grant from the UK Medical Research Council (MR/M021475/1). RP is supported by a Medical Research Council Professorship award (G19/2). T.A. McAdams and Y.I. Ahmadzadeh are funded by a Sir Henry Dale Fellowship awarded to T.A. McAdams, jointly funded by the Wellcome Trust and the Royal Society (grant number 107706/Z/15/Z). R. Cheesman is supported by an ESRC studentship. We thank Aysu Okbay and the SSGAC for providing genome-wide association summary statistics for educational attainment excluding UK Biobank.

Additional information

Supplemental Material is available for this paper.

Author contributions

R.C. and G.B conceived and designed the study. R.C. analysed the data and wrote the manuscript. Genotype data quality control was conducted by A.H. and J.C. All co-authors provided critical revisions and approved the final version of the manuscript for submission.

References

- Allen, N. E., Sudlow, C., Peakman, T., Collins, R., & UK Biobank. (2014). UK biobank data: come and get it. *Science Translational Medicine*, 6(224), 224ed4. doi:10.1126/scitranslmed.3008601
- Bates, T. C., Maher, B. S., Medland, S. E., McAloney, K., Wright, M. J., Hansell, N. K., ... Gillespie, N. A. (2018). The Nature of Nurture: Using a Virtual-Parent Design to Test Parenting Effects on Children's Educational Attainment in Genotyped Families. *Twin Research and Human Genetics*, 21(2), 73–83. doi:10.1017/thg.2018.11
- Brown, M. B., & Forsythe, A. B. (1974). Robust tests for the equality of variances. *Journal of the American Statistical Association*, 69(346), 364–367. doi:10.1080/01621459.1974.10482955
- Bulik-Sullivan, B., Finucane, H. K., Anttila, V., Gusev, A., Day, F. R., Loh, P.-R., ... Neale. (2015). An atlas of genetic correlations across human diseases and traits. *Nature Genetics*, 47(11), 1236–1241. doi:10.1038/ng.3406
- Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Schizophrenia Working Group of the Psychiatric Genomics Consortium, ... Neale, B. M. (2015). LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature Genetics*, 47(3), 291–295. doi:10.1038/ng.3211
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., ... Marchini, J. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726), 203–209. doi:10.1038/s41586-018-0579-z
- Choi Shing Wan. (n.d.). PRSice 2 Wiki · GitHub. Retrieved January 31, 2018, from <https://github.com/choishingwan/PRSice/wiki>
- Domingue, B. W., Belsky, D., Conley, D., Harris, K. M., & Boardman, J. D. (2015). Polygenic Influence on Educational Attainment: New evidence from The National Longitudinal Study of Adolescent to Adult Health. *AERA Open*, 1(3), 1–13. doi:10.1177/2332858415599972
- Domingue, B. and Fletcher, J. (2019). Separating measured genetic and environmental effects: evidence linking parental genotype and adopted child outcomes. BioRxiv.
- Duyme, M., Dumaret, A.C. and Tomkiewicz, S. (1999). How can we boost IQs of “dull children”? A late adoption study. *Proceedings of the National Academy of Sciences of the United States of America* 96(15), pp. 8790–8794.
- Evans, L. M., & Keller, M. C. (2018). Using partitioned heritability methods to explore genetic architecture. *Nature Reviews. Genetics*, 19(3), 185. doi:10.1038/nrg.2018.6
- Fry, A., Littlejohns, T.J., Sudlow, C., et al. (2017). Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank Participants With Those of the General Population. *American Journal of Epidemiology* 186(9), pp. 1026–1034.
- Holmlund, H., Lindahl, M., & Plug, E. (2011). The causal effect of parents' schooling on children's schooling: A comparison of estimation methods. *Journal of Economic Literature*, 49(3), 615–651. doi:10.1257/jel.49.3.615
- Howard, J.A., Smith, S.L. and Ryan, S.D. (2004). A Comparative Study of Child Welfare Adoptions with Other Types of Adopted Children and Birth Children. *Adoption quarterly* 7(3), pp. 1–30.
- Keller, M. C. (2014). Gene × environment interaction studies have not properly controlled for potential confounders: the problem and the (simple) solution. *Biological Psychiatry*, 75(1), 18–24. doi:10.1016/j.biopsych.2013.09.006
- Keyes, K. M., & Westreich, D. (2019). UK Biobank, big data, and the consequences of non-representativeness. *The Lancet*, 393(10178), 1297. doi:10.1016/S0140-6736(18)33067-8
- Kong, A., Thorleifsson, G., Frigge, M. L., Vilhjalmsón, B. J., Young, A. I., Thorgeirsson, T. E., ... Stefansson, K. (2018). The nature of nurture: Effects of parental genotypes. *Science*, 359(6374), 424–428. doi:10.1126/science.aan6877
- Krapohl, E., Hannigan, L. J., Pingault, J. B., Patel, H., Kadeva, N., Curtis, C., ... Plomin, R. (2017). Widespread covariation of early environmental exposures and trait-associated polygenic variation. *Proceedings of the National Academy of Sciences of the United States of America*, 114(44), 11727–11732. doi:10.1073/pnas.1707178114

- Lee, J. J., Wedow, R., Okbay, A., Kong, E., Maghzian, O., Zacher, M., ... Karlsson Linnér, R. (2018). Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nature Genetics*, *50*(8), 1112–1121. doi:10.1038/s41588-018-0147-3
- Loehlin & De Fries. (1987). Genotype-environment correlation revisited. *Behavior Genetics*, *(22)*, 731–732.
- Maughan, B., Collishaw, S., & Pickles, A. (1998). School achievement and adult qualifications among adoptees: a longitudinal study. *Journal of Child Psychology and Psychiatry, and Allied Disciplines*, *39*(5), 669–685. doi:10.1111/1469-7610.00367
- Natsuaki, M. N., Neiderhiser, J. M., Harold, G. T., Shaw, D. S., Reiss, D., & Leve, L. D. (2019). Siblings reared apart: A sibling comparison study on rearing environment differences. *Developmental Psychology*, *55*(6), 1182–1190. doi:10.1037/dev0000710
- Plomin, R., DeFries, J. C., & Loehlin, J. C. (1977). Genotype-environment interaction and correlation in the analysis of human behavior. *Psychological Bulletin*, *84*(2), 309–322. doi:10.1037/0033-2909.84.2.309
- Plomin, Robert. (1994). *Genetics and Experience: the Interplay Between Nature and Nurture*. SAGE Publications.
- Plomin, Robert. (2014). Genotype-environment correlation in the era of DNA. *Behavior Genetics*, *44*(6), 629–638. doi:10.1007/s10519-014-9673-7
- Plomin, R., DeFries, J.C., Knopik, V.S. and Neiderhiser, J.M. (2016). Top 10 replicated findings from behavioral genetics. *Perspectives on psychological science : a journal of the Association for Psychological Science* *11*(1), pp. 3–23.
- Plomin, Robert, Loehlin, J. C., & DeFries, J. C. (1985). Genetic and environmental components of “environmental” influences. *Developmental Psychology*, *21*(3), 391–402. doi:10.1037/0012-1649.21.3.391
- Plomin, Robert, & von Stumm, S. (2018). The new genetics of intelligence. *Nature Reviews. Genetics*, *19*(3), 148–159. doi:10.1038/nrg.2017.104
- Rutter, M. (2006). *Genes and Behavior: Nature-Nurture Interplay Explained*. Wiley.
- Selzam, S., Ritchie, S. J., Pingault, J.-B., Reynolds, C. A., O’Reilly, P. F., & Plomin, R. (2019). Comparing within- and between-family polygenic score prediction. *BioRxiv*. doi:10.1101/605006
- Turkheimer, E. (1991). Individual and group differences in adoption studies of IQ. *Psychological Bulletin*, *110*, 392–405.
- van Ijzendoorn, M. H., Juffer, F., & Poelhuis, C. W. K. (2005). Adoption and cognitive development: a meta-analytic comparison of adopted and nonadopted children’s IQ and school performance. *Psychological Bulletin*, *131*(2), 301–316. doi:10.1037/0033-2909.131.2.301
- Wertz, J., Belsky, J., Moffitt, T. E., Belsky, D. W., Harrington, H., Avinun, R., ... Caspi, A. (2019). Genetics of nurture: A test of the hypothesis that parents’ genetics predict their observed caregiving. *Developmental Psychology*. doi:10.1037/dev0000709
- Wertz, J., Moffitt, T. E., Agnew-Blais, J., Arseneault, L., Belsky, D. W., Corcoran, D. L., ... Caspi, A. (2018). Using DNA from mothers and children to study parental investment in children’s educational attainment. *BioRxiv*. doi:10.1101/489781
- Willer, C. J., Li, Y., & Abecasis, G. R. (2010). METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics*, *26*(17), 2190–2191. doi:10.1093/bioinformatics/btq340
- Wood, A. R., Esko, T., Yang, J., Vedantam, S., Pers, T. H., Gustafsson, S., ... Kutalik, Z. (2014). Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature Genetics*, *46*(11), 1173–1186. doi:10.1038/ng.3097
- Yang, J., Lee, S. H., Goddard, M. E., & Visscher, P. M. (2011). GCTA: a tool for genome-wide complex trait analysis. *American Journal of Human Genetics*, *88*(1), 76–82. doi:10.1016/j.ajhg.2010.11.011
- Young, A. I., Frigge, M. L., Gudbjartsson, D. F., Thorleifsson, G., Bjornsdottir, G., Sulem, P., ... Kong, A. (2018). Relatedness disequilibrium regression estimates heritability without environmental bias. *Nature Genetics*, *50*(9), 1304–1310. doi:10.1038/s41588-018-0178-9

Zheng, J., Erzurumluoglu, A. M., Elsworth, B. L., Kemp, J. P., Howe, L., Haycock, P. C., ... Neale, B. M. (2017). LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics*, 33(2), 272–279. doi:10.1093/bioinformatics/btw613