**MULTI-VARIABLE APPROACHES TO POLYGENIC TRAIT PREDICTION**

Krapohl, Eva Maria Laura

*Awarding institution:*
King's College London

**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.

# MULTI-VARIABLE APPROACHES TO POLYGENIC TRAIT PREDICTION

Eva Krapohl

MRC Social, Genetic and Developmental Psychiatry Centre

Institute of Psychiatry, Psychology and Neuroscience

King's College London

Submitted for the degree of Doctor Of Philosophy in Statistical Genetics

2017

# Author declaration

The research presented in this thesis was conducted as part of the Twin Early Development Study (TEDS), an ongoing longitudinal study following twins born in England and Wales between 1994 and 1996. Phenotypic and genomic data were collected and cleaned by people other than myself. For chapters 3 and 5, I conducted the quality control, imputation, and harmonisation of the genetic data, and I prepared both phenotypic and genomic data for all the analyses presented this thesis. In all other respects, to the best of my knowledge, the work presented in this thesis is original and my own work, except where acknowledged in the text.

# Acknowledgements

# Abstract

Robust evidence for the polygenicity and genetic correlations of complex traits across the phenome suggests both the necessity of polygenic instruments and the value of multi-trait prediction models. This thesis used multi-variable approaches in four papers and along two main threads:

**Multi-variable approaches to trait prediction**  A primary goal of polygenic scores, which aggregate effects of trait-associated variants discovered in genome-wide association studies (GWAS), is to estimate individual-specific genetic propensities to predict trait variation. This is typically achieved using one polygenic score predicting one outcome. Extending this to a multi-variable approach, a 'phenome-wide analysis of genome-wide polygenic scores' mapped associations between 13 polygenic scores created from GWAS for psychiatric disorders and cognitive traits and 50 behavioural traits.

Extending the multi-variable approach further, a multi-polygenic score approach was employed to increase prediction by exploiting the joint power of multiple discovery GWAS in the same model. A regularised regression model combining summary statistics of 81 trait GWAS improved out-of-sample prediction of three child outcomes over the best single-predictor model.

**Multi-variable approaches to gene-environment correlation**  Although gene-environment correlation is widely investigated by family studies and recently by SNP-heritability studies, the possibility that genetic effects on traits capture environmental risk factors or protective factors has been neglected by polygenic prediction models. First, a study using genome-wide SNP-heritability estimation and polygenic score analysis provided the first molecular evidence for substantial genetic influence on differences in children's educational achievement and its association with family socio-economic status.

Second, covariation between offspring trait-associated polygenic variation and a wide range of parent-mediated environmental exposures was estimated. For this, a mixed linear model estimated the effects of multiple polygenic scores on each environmental exposure while controlling for overall relatedness by fitting the effects of all SNPs as random effects. Findings illustrate the relevance of gene-environment correlation for polygenic prediction models.

Taken together, the analyses illustrate the value of multi-variable approaches to complex trait prediction, as well as their current limitations and future potential.

# Contents

# List of Figures

## 3 Phenome-wide analysis of genome-wide polygenic scores

## 4 Multi-polygenic score approach to trait prediction

## 5 Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs

## 6 Widespread covariation of early environmental exposures and trait-associated polygenic variation

# 1 General Introduction

Evidence across multiple methods converges in showing ubiquitous and substantial genetic influence on human trait variation (Ge et al., 2017; Muñoz et al., 2016; Polderman et al., 2015). Yet, success in making predictions of complex trait phenotypes from genotype data has been complicated by two main challenges: genetic effects on trait variation are spread across genetic loci and they are shared among multiple traits. This suggests both the necessity of polygenic instruments and the value of multi-variable models for trait prediction. A further complication arises from the empirical observation that genetic correlations between traits extend to the environment: genetic and environmental variation is not independent; individuals' environmental exposure partially depends on their genotype. This thesis used polygenic multi-variable approaches for trait prediction and for investigating covariance between traits and environments.

The following sections provide an outline of the phenomena of polygenicity and pleiotropy, and how multi-trait approaches can leverage them for genomic prediction.

## 1.1 Polygenicity

The investigation of the genetic basis of human traits and diseases is made difficult by the empirical observation that common traits are complex (Robinson et al., 2014). That is, trait variation is a function of a multitude of genetic and environmental factors and their interplay. Complex traits stand in contrast to Mendelian traits, which are strongly influenced by variation within a single gene and are characterised by their classic patterns of inheritance within families (Garrod, 1902; Mendel, 1866). Traits that strictly conform to Mendelian principles of inheritance are relatively rare. Most traits do not follow readily predictable patterns of inheritance, but are subject to a wide spectrum of non-Mendelian, multifactorial mechanisms. Often the genetic variants affecting complex traits are large in quantity but small in effect (Falconer and Mackay, 1996).

In 1918, Fisher seminally proposed that continuous variation amongst phenotypic traits can be the result of Mendelian inheritance if multiple genetic variants affect the trait (Fisher, 1918). Due to the multitude of possible allelic combinations, traits under polygenic inheritance follow a normal continuous distribution. As the num-

ber of genetic variants approaches infinity, individual effects approach zero (Barton, 1990). Fisher's polygenic model conceptually reconciled the quantitative and qualitative perspectives of genetic influence on trait variation.

In the last decades, the advent and development of microarrays that can measure millions of genetic markers have led to an explosion of genome-wide association studies (GWAS) of complex traits and disorders. This has empirically reinforced the polygenic trait model, with the extent of polygenicity being even more extreme than anticipated.

For many human traits, single-nucleotide polymorphism (SNP)-trait associations have been identified through GWAS. To identify trait-associated loci, a large series of simple linear variant-trait regressions is conducted across the genome, each testing whether allele frequency at a given locus is significantly associated with phenotypic variation. In response to the ubiquitous polygenicity of complex traits, ever-larger GWAS are being conducted to increase statistical power for detecting variants with small effect sizes.

Although genome-wide association studies have been successful in discovering and replicating SNP associations for many traits and disorders (Duncan et al., 2017; Robinson et al., 2014; Visscher et al., 2012, 2017; Zheng et al., 2017), the lack of larger SNP-trait associations in well-powered GWAS provides evidence for the ubiquitous heritability of complex dimensions and common disorders being spread across tens or hundreds of thousands of common DNA variants with individually tiny effect. Evidence for the heritability contributed by each chromosome being proportional to its physical lengths suggests that effects might be relatively evenly spread across the genome (Shi et al., 2016a; Visscher et al., 2006).

## 1.2  Genotype-based trait prediction

The polygenicity of complex traits poses an immense conceptual and statistical challenge for deciphering the links between genetic and phenotypic variation, and therefore complicating genotype-based trait prediction.

Next to discovery of trait-associated variants and their biological function, there is increasing interest in predicting trait variation from genotype data, which is the focus of this thesis. These predictions rely on the estimation of the effects of genetic variants in a discovery sample, with subsequent validation in an independent sample,

$$R^2 = \frac{h^2_{M_{eff}}}{1 + \frac{M_{eff}}{N_{eff}h^2_{M_{eff}}}(1 - R^2)}$$

Equation 1.1: Proportion of phenotypic variance explained by a predictor of a quantitative trait

and eventually prediction of individuals' future phenotypes in practice. In contrast to the deterministic genetic tests for fully penetrant Mendelian disorders, genetic predictions for complex traits are probabilistic.

Regardless of genetic architecture of the predicted trait, an approximation of the predictive power of a genotype-based predictor using the estimated effects of all markers is a function of the effective number of independently associated variants ($M_{eff}$), the trait variance they account for ($h^2_{M_{eff}}$), and the effective sample size of the discovery sample ($N_{eff}$), equation 1.1 (Daetwyler et al., 2008; Dudbridge, 2013; Goddard, 2009; Palla and Dudbridge, 2015; Visscher et al., 2010; Wray et al., 2014, 2013).

$R^2$ is the accuracy of the genetic predictor, with $R$ being the correlation between the predictor and the to be predicted outcome trait. The efficacy of a genetic predictor can be approximated by its predictive power relative to the best possible predictor or the overall heritability; $\frac{R^2}{h^2}$. Discrepancy between $R^2$ and $h^2$ is therefore affected by the parameters in equation 1.1.

### 1.2.1 Current limitations of genotype-based prediction

Making predictions of complex trait phenotypes from genotype data has several conceptual and logistical limitations. Genomic prediction is limited by the overall heritability ($h^2$), that is, the proportion of trait variation explained by genetic variation. Hence, a linear coefficient of determination ($R^2$) relying solely on genetic variants is unable to explain all trait variation in complex traits where $h^2$ is less than 100%. What is more, the condition for $R^2$ being equal to $h^2$ is that all genetic effects are perfectly reliably estimated.

The next paragraphs describe current limitations of genotype-based polygenic trait prediction, followed by an outline of how the first part of this thesis was aimed at exploring the potential of increasing $R^2$ for polygenic traits using a multi-variable approach.

Known pedigree structure has been an important methodological resource for quantifying genetic influences on trait variation. Pedigree-based heritability estimates, $h^2_{pedigree}$, are derived by fitting the covariance structure specified by the matrix of known kinship coefficients to a vector of measured phenotypes. Heritability estimates from family studies indicate ubiquitous additive genetic influence on complex traits (Polderman et al., 2015), suggesting great potential for genotype-based trait prediction. So far, genotype-based complex trait prediction has not reached this upper limit for prediction. Multiple possible sources of this discrepancy have been proposed (Eichler et al., 2010; Manolio et al., 2009).

Genotype-based trait prediction is limited by $h^2_M$ (see Eq. 1.1). For almost all complex traits, $h^2_M$, estimated as $h^2_{SNP}$ by fitting the additive effects of all measured SNPs as random effects in a linear mixed model (Speed et al., 2012; Yang et al., 2011a; Zhou and Stephens, 2014), is systematically smaller than additive $h^2_{pedigree}$ estimated by pedigree-based methods (Liu et al., 2015; Muñoz et al., 2016; Yang et al., 2015; Zaitlen et al., 2013, 2014).

Model misspecifications in estimating $h^2_{SNP}$ or $h^2_{pedigree}$ might be one source of this discrepancy. First, family-based approaches such as twin studies might be biased upwards by epistatic interactions, or more likely by shared environment, with improved parameterisation of shared familial environmental effects showing reduced overestimation (Liu et al., 2015; Muñoz et al., 2016; Yang et al., 2015; Zaitlen et al., 2014, 2013). Second, it has been shown that conventional $h^2_{SNP}$ estimation models might underestimate $h^2_{SNP}$ due to incorrect prior assumptions about the distribution of genetic effects across the genome as a function of minor allele frequency and linkage disequilibrium (LD) (Speed et al., 2017). These two lines of evidence suggest that the discrepancy between $h^2_M$ and $h^2_{pedigree}$ may be smaller than previously thought.

$h^2_M$ only reflects variants correlated with the common markers included on SNP arrays, i.e. $h^2_{SNP}$. Therefore, $h^2_{SNP}$ underestimates additive genetic variance due to imperfect LD between causal and genotyped (or imputed) SNPs. Common SNP arrays take advantage of the fact that around half a million common SNPs in the human genome 'tag' most of the common variation (in non-African populations) (The International HapMap Consortium, 2005). Although common ($>5\%$) conventionally genotyped (and imputed) markers capture variation of unmeasured common markers very well, they do not reliably 'tag' low frequency alleles. That is because LD of two loci strongly depends on their allele frequency, with low frequency alleles being weakly 'tagged' by proximal genotyped common alleles (Wray, 2005). By design models using SNP from common 'SNP chips' are therefore limited to identify associations with causal variants that are relatively frequent in the population. Variants conferring large disadvantageous effects are screened by natural selection and held

at low population frequencies. Because SNP chips are limited to the high allelic frequency spectrum, models that rely on these SNPs are targeting variants with small effects only.

Due to the limited power of GWAS in face of the polygenic architecture of most complex traits and the technological limitation to common variants, genome-wide significant variants discovered by GWAS explain only a small fraction of the additive genetic variance tagged by all investigated markers, $h^2_{GWAS} < h^2_{SNP}$. This discrepancy has been termed 'hidden heritability': the cumulative trait variation explained by GWAS 'hits' lags behind their known 'net effect' estimated by $h^2_{SNP}$ methods, because small effect variants are 'hiding' within the noise below the genome-wide significance threshold due to insufficient statistical power to detect them. For example, using SNP-heritability estimation approaches, it has been estimated that over half the variation in human height can be attributed to the common SNPs on a genome-wide genotyping array. In comparison, only ∼16% phenotypic variation can be explained by ∼670 SNPs reaching genome-wide significance (Wood et al., 2014; Yang et al., 2010, 2011b).

### 1.2.2 Polygenic prediction approaches

Polygenic prediction methods have been used with the aim to increase $R^2$ from $h^2_{GWAS}$ towards the current ceiling of genotype-based prediction, i.e. $h^2_{SNP}$. For the same set of SNPs, $R^2$ would be equal to $h^2_M$ only if the marker effects were estimated without error by the discovery GWAS. Therefore, the aim of polygenic prediction methods is to optimise the balance between the predictive gain from including more SNPs into the model and the predictive loss from the markers‚Äô unreliably estimated effects in the discovery sample. By including variants below the threshold for genome-wide significance, polygenic scores can reveal some of the "hidden heritability", i.e. the gap in trait variation explained by genome-wide significant SNPs discovered by GWAS and variation tagged by all common SNPs together.

Polygenic scores, which aggregate the effects of thousands of trait-associated genetic variants discovered in GWAS, have been widely applied to predict trait variation in independent target samples (Evans et al., 2009; Marioni et al., 2016; Power et al., 2015; Purcell et al., 2009; Szulkin et al., 2015; Vassos et al., 2017; Vassy et al., 2014; Wray et al., 2007). However, polygenic scores have typically only explained a fraction of the trait variance of polygenic traits, serving as validation tools for SNP-trait

associations discovered by GWAS rather than predictive tools. Major advances in the prediction approaches are required for genotype-based polygenic trait prediction to become relevant in practice by achieving meaningful risk stratification.

Different lines of approaches have been pursued with the aim of increasing predictive power of polygenic score prediction. The main strategy has been to improve statistical power of detecting associations between SNPs and polygenic traits by increasing the sample of the discovery set (Visscher et al., 2017). Less focus has been on increasing prediction accuracy by improving models used when predicting into the target sample. Here strategies have included improving modelling of LD (Vilhjalmsson et al., 2015) or optimising number of SNPs included in the prediction model (Euesden et al., 2015). The first half of the current work explored the potential of increasing trait prediction by leveraging genetic signal captured by multiple discovery sets.

The following section provides an overview of evidence for genetic correlations across the phenome and the potential of increasing trait prediction by leveraging the genetic correlations across traits in polygenic multi-trait prediction models.


## 1.3 Genetic correlation and multi-trait approaches

A fully 'infinitesimal' polygenic trait model would necessarily imply shared genetic loci between traits to the extent that traits are under genetic influence. However, even more modest extents of polygenicity makes it likely that some of the genetic signal (positively or negatively) correlates between traits (Wright, 1984). This is consistent with empirical data.

Next to polygenicity, the other defining characteristic of complex traits and common disorders is the ubiquity of genetic correlations between them. Genetic correlation is the empirical observation that genetic variation associated with one trait co-varies with that associated with another trait. Mendelian mutations resulting in specific syndromes or diseases are often associated with multiple phenotypes in an affected individual. For polygenic traits, there exists robust evidence across twin and molecular methods for genetic correlations between psychiatric disorders, between anthropometric traits, and between educational and cognitive traits, as well as for genetic correlations across these realms (Bulik-Sullivan et al., 2015a; Calvin et al., 2012; PGC, 2013; Davis et al., 2014; Duncan et al., 2017; Kovas and Plomin, 2006; Krapohl et al., 2014; Lichtenstein et al., 2009; Pickrell et al., 2016; Visscher and Yang, 2016).

Multivariate approaches have long been used by pedigree designs to investigate shared genetic aetiology (Kendler et al., 2008; Krapohl et al., 2014; Lichtenstein et al., 2009; Plomin and DeFries, 1979). More recently, molecular genetic approaches have been developed to investigate genetic correlation between two traits, either using individual-level genotype or GWAS summary statistics (Bulik-Sullivan et al., 2015a; Johnson, 2013; Lee et al., 2012; Zhou and Stephens, 2014). Cross-trait polygenic score analyses have replicated and identified shared genetic aetiology between traits, such as predicting addiction from genetic liability for schizophrenia and bipolar or cognitive ability from genetic variation associated with neural structures (Luciano et al., 2015; Reginsson et al., 2017).

Researchers have also started to be interested in simultaneously analysing multiple correlated traits to improve statistical power of the discovery data by leveraging cross-trait covariance (Baselmans et al., 2017; Bolormaa et al., 2014; Ferreira and Purcell, 2009; Korte et al., 2012; Maier et al., 2015; Rietveld et al., 2014; Shim et al., 2015; Turley et al., 2017). These approaches rely on substantial and consistent correlations between discovery GWAS. Some of these approaches are designed to discriminate heterogeneity and homogeneity in SNP-trait associations (Bhattacharjee et al., 2012; Majumdar et al., 2017) or different types of direct and indirect pleiotropy (Giambartolomei et al., 2014; Pickrell et al., 2016; Shi et al., 2016b; Stephens, 2013) across correlated traits.

### 1.3.1  Overview of chapters 3 and 4: Multi-variable approaches to trait prediction

A primary goal of polygenic scores is to estimate individual-specific genetic propensities. This is typically achieved using a single polygenic score to predict a single outcome trait.

Chapter 3 describes a systematic investigation of profiles of associations between multiple genome-wide polygenic scores across a wide range of behavioural traits. Specifically, this 'phenome-wide analysis of genome-wide polygenic scores' mapped associations between 13 polygenic scores created from GWAS for psychiatric disorders and cognitive traits and 50 behavioural traits measured in adolescence.

When the aim is prediction, genetic correlations between traits can be used for maximising prediction power within a multi-variable approach, regardless of the underlying mechanisms. Therefore, the premise of the approach used in chapter 4 was to maximise prediction of developmental outcomes, rather than investigating their

aetiology. This stands in contrast to multi-trait meta-analytic approaches of GWAS summary statistics, which rely on substantial and consistent correlations between discovery GWAS and whose main aim is variant discovery (Baselmans et al., 2017; Ferreira and Purcell, 2009; Maier et al., 2015; Turley et al., 2017) The multi-polygenic score approach used here allowed for, but does not require, correlation among polygenic predictors.

Chapter 4 employed a multi-polygenic score approach to increase predictive power by exploiting the joint power of multiple discovery GWAS in the same model, without assumptions about the relationships among predictors. I selected GWAS from a centralised repository of summary statistics – based on their statistical power and regardless of prior evidence for association with the outcomes – to predict three core developmental outcomes in our independent target sample: educational achievement, body-mass index, and general cognitive ability. Using repeated cross-validation, I trained and validated the prediction models using elastic net regularised regression. Finally, I compared out-of-sample prediction of these multi-score models to the best single predictor models.

## 1.4    Genotype-environment correlation

The heritability, polygenicity, and genetic correlations observed within the realm of phenotypes also cross over to that of environments; genetic variation and variation in enviornmental exposures is not independent.

Converging evidence from family, twin, and adoption studies has shown that individuals' exposure to environments and perceptions of environments varies as a function of their genotype. This genotype-environment correlation includes both parenting characteristics and broad socio-economic variables (Avinun and Knafo, 2013; Butcher and Plomin, 2008; Kendler and Baker, 2007; Klahr and Burt, 2014; Plomin and Bergeman, 1991; Vinkhuyzen et al., 2010). In the past decade, quantitative genetic research has started investigating genetic and environmental contributions to correlations between environmental factors and children's developmental outcomes (Colen and Ramey, 2014; D'Onofrio et al., 2010a,b, 2007; Evenhouse and Reilly, 2005; Larsson et al., 2014).

Some newer designs such as the children-of-twins designs allow for disentangling different types of genotype-environment correlation and identify environmental influences controlling for genetic confounds (Harden et al., 2007; Knopik et al., 2006;

Lynch et al., 2006; Narusyte et al., 2008; Silberg et al., 2010). However, these designs are limited by the extent to which environmental variables differ between close relatives.

Corroborating evidence for gene-environment correlation comes from $h^2_{SNP}$ studies that estimate the 'net' genetic effect on trait variation from empirical genomic similarity in unrelated individuals. Using $h^2_{SNP}$ estimation, studies have shown variation in individuals' social deprivation, household income, stressful life events, and family socio-economic status partially reflects individuals' differences across genome-wide common genetic variants measured on SNP arrays (Benjamin et al., 2012; Davies et al., 2015; Hill et al., 2016; Marioni et al., 2014; Power et al., 2013; Trzaskowski et al., 2014). There have also been a few reports of extending SNP heritability analysis to estimate genetic correlations between environmental measures and measures of children's developmental outcomes (Davies et al., 2015; Trzaskowski et al., 2014).

Gene-environment correlation is a common subject of investigation by family studies and recently by SNP-heritability studies. However, the possibility that individuals' trait-associated polygenic variation captures variance in established environmental risk and protective factors has not been considered by polygenic trait prediction models, which use genetic variants identified by GWAS to estimate individual-specific genetic trait propensities.

### 1.4.1 Overview of chapters 5 and 6: Multi-variable approaches to gene-environment correlation

Early environmental exposures are amongst the best predictors for health and educational outcomes. For instance, maternal smoking during pregnancy, watching television, harsh parenting, and older paternal age have been identified as risk factors for a range of behaviour problems, whereas higher parental socio-economic status and breastfeeding have been associated with more favourable child outcomes (Afifi et al., 2012; Ainsworth, 2002; Bender et al., 2007; Byrne et al., 2003; Caspi et al., 2016; Danner, 2008; de Kluiver et al., 2017; D'Onofrio et al., 2014; Eamon, 2005; Garner and Raudenbush, 1991; Gentile et al., 2004; Gershoff, 2002; Huizink and Mulder, 2006; Jago et al., 2005; Janecka et al., 2017; Knox, 2010; Leventhal and Brooks-Gunn, 2000; Malaspina, 2001; Räsänen et al., 1999; Reichenberg et al., 2006; Sandin et al., 2016; Sirin, 2005; Taylor et al., 2010; Victora et al., 2015; White et al., 1999; White, 1982).

Chapter 5, using genome-wide SNP-heritability estimation and polygenic score ana-

lysis, provided the first molecular evidence for substantial genetic influence on differences in children's educational achievement and its association with family socio-economic status (SES). The analyses also tested to what extent the observed genetic covariation between children's educational achievement and family SES was explained by children's general cognitive ability.

Chapter 6 investigated to what extent offspring trait-associated alleles covary with parental traits and behaviours previously reported to be environmental risk or protective factors for important child outcomes. Specifically, mixed linear models estimated the effects of trait-associated polygenic variation while controlling for overall genetic relatedness by fitting the effects of all SNPs as random effects. A second set of analyses tested to what extent offspring genetic trait propensities contribute to the correlation between parenting characteristics and children's developmental outcomes.

## 1.5 Summary

Robust evidence for the polygenicity and genetic correlations of complex traits across the phenome and environment suggests both the necessity of polygenic instruments and the value of multi-trait models. This thesis uses several multi-variable genome-wide approaches for trait prediction (chapters 3 and 4), and to investigate genotype-environment correlation relevant for polygenic prediction (chapters 5 and 6).

The next chapter provides an overview of some of the methods used in these studies, methodological details are described in the later chapters.

# 2    General methods

While all methods employed in this thesis are explained in detail within the respective chapters, the following sections provide a basic outline of the two main techniques used: estimation of genetic variance and covariance from genome-wide SNPs (chapter 3) and polygenic score prediction (chapters 3, 4, 5, and 6).

## 2.1    Estimation of genetic variance and covariance from genome-wide SNPs measured in unrelated individuals

All quantitative genetics models estimate genetic and residual contribution to trait variation by fitting a covariance structure specified by a matrix of kinship coefficients to a vector of measured phenotypes (Falconer and Mackay, 1996; Wright, 1920). Models vary however in which kind of genetic variation they consider and how it is measured.

The generic model estimating the extent to which phenotypic similarity between pairs of individuals is accounted for by their genetic relatedness can be expressed in terms of variance components:

$$Var(Y) = \sigma_a^2 A + \sigma_d^2 D + \sigma_c^2 C + \sigma^2 I$$

Equation 2.1: Generic variance component model

With $Y$ being a vector of phenotype values, $A$ and $D$ matrices of kinship coefficients corresponding to additive and dominant genetic effects, $C$ common environmental effects, and $I$ is an identify matrix containing residual ('non-shared environmental') effects.

The widespread availability of high-throughput SNP-chip genotyping has enabled the move from using known kinship coefficients from pedigree data to estimating genomic differences from high-density SNP arrays (Meuwissen et al., 2001; Mousseau et al., 1998; Ritland, 1996; Thomas et al., 2002; van Kleunen and Ritland, 2005, 2004; Visscher et al., 2007). Previously used in animal breeding and plant genetics, recently this has attracted interest and methodological development in human genetics because it offers several advantages over traditional pedigree-based methods (Yang et al., 2013, 2011a).

Whereas traditional pedigree approaches use known kinship coefficients, $h_{SNP}^2$ estim-

ation methods fit empirically established genomic similarity between conventionally unrelated individuals (Meuwissen et al., 2001; Speed et al., 2012; Yang et al., 2011a). Here, genetic similarity is estimated from short segments of nucleotide sequences shared between unrelated individuals. Importantly for this thesis, the reliance on empirical genomic similarity in 'unrelated' individuals allows for decomposition of phenotypic variance/covariance of family-level variables.

$h^2_{SNP}$ methods estimate additive genetic variance captured by genome-wide SNPs by fitting the effects of all sampled (i.e. genotyped or imputed) SNPs as random effects in a linear mixed model. Linear mixed models, ubiquitously used in heritability estimation (Henderson et al., 1959; Robinson, 1991), fit a covariance structure specified by a matrix of kinship coefficients to a vector of measured phenotypes. "Mixed" models contain both an unobserved random effect, usually interpreted in terms of a polygenic contribution to the trait, and fixed effects (i.e. 'covariates', not shown in the equations below). $h^2_{SNP}$ is estimated by the squared regression coefficient of the random effect. In relation to the generic model (Eq. 2.1), the additive genetic component representing $h^2_{SNP}$ can be expressed as:

$$h^2_{SNP} = \frac{\sigma^2_a}{\sigma^2_a + \sigma^2_d + \sigma^2_c + \sigma^2_i}$$

Equation 2.2: $h^2_{SNP}$ 'SNP-heritability'

Using restricted/residual maximum likelihood (REML), the $h^2_{SNP}$ estimation model partitions trait variance into additive polygenetic ($a$) and residual effects ($i$), which also contains all non-additive genetic effects.

$$Var(Y) = \sigma^2_a A + \sigma^2_i I$$

Equation 2.3: Univariate mixed linear model for $h^2_{SNP}$ estimation

The genetic component of the decomposition is a product of genomic similarity ($A$) and the estimated additive genetic effect ($\sigma^2_a$). The genomic similarity matrix ($A$) contains pairwise genomic similarity between all pairs of individuals in the samples. The genomic similarity between a given pair of individuals is their genome-wide allelic correlation, weighted by allele frequencies for each SNP.

$$\frac{1}{M} \sum_{i=1}^{M} \frac{(G_{ij} - 2p_i)(G_{ik} - 2p_i)}{2p_i(1 - p_i)}$$

Equation 2.4: Algorithm for the estimation of pairwise genetic similarity (Yang et al., 2011a)

$G_{ij}$ is individual $j$'s genotype (i.e. number of copies of the reference allele) at $SNP_i$, with $p_i$ being the frequency of the reference allele. $M$ is the number of markers in the model. Notably, causal variants only contribute if they are tagged by the measured SNPs, creating a bias towards common causal variants, which are better tagged than rare causal variants.

This model can be extended to the bivariate (or multivariate) level by relating the pairwise genetic similarity matrix to a phenotypic covariance matrix between traits (Lee et al., 2012; Thompson, 1973; Zhou and Stephens, 2014).

$$(V) = \left[ \begin{array}{cc} Z_1 A Z_1' \sigma_{a1}^2 + I \sigma_{i1}^2 & Z_1 A Z_2' \sigma_{a1,a2}^2 \\ Z_2 A Z_1' \sigma_{a1,a2}^2 & Z_2 A Z_2' \sigma_{a2}^2 + I \sigma_{i2}^2 \end{array} \right]$$

Equation 2.5: Variance/covariance matrix of bivariate mixed linear model (Yang et al., 2011a)

$A$ is the genetic similarity matrix, and $Z$ is an incidence matrix holding the random genetic effects. Here $\sigma_{a1,a2}^2$ estimates the genetic covariance between the two traits, with the residual component modelled only on the variance of each trait. Because the genetic correlation estimate is a non-linear function of the genetic variances and covariances of the traits, there is no explicit way to estimate its variance. In the commonly used Genome-wide Complex Trait Analysis (GCTA) software, the sampling distribution of the statistic is calculated using an average information method to obtain the standard error of each estimate (Lee et al., 2012; Lee and van der Werf, 2006).

Genetic correlation is the ratio between genetic covariance and the genetic variances of the two traits, all of which are subject to the same underestimation. Therefore, the estimate of genetic correlation will be similar to that of the twin method, if the degree of correlation is the same for the co/variance tagged as for the co/variance untagged by the sampled SNPs.

## 2.2 Polygenic score prediction

Because the individual effects of common genetic variants on complex traits are miniscule, single-variant models are of little use for prediction. When the aim is prediction, polygenic scores can be used to aggregate the effects of multiple genetic variants discovered by independent GWAS (Dudbridge, 2013; Palla and Dudbridge, 2015; Purcell et al., 2009). Rather than just aggregating SNPs passing the level of genome-wide significance, a recent development is to aggregate a much larger number of SNPs, weighted by their GWAS effect size estimate. Unlike quantitative genetic

designs that estimate the net effect captured by the measured genetic differences in the population, polygenic scores provide individual-specific estimates of genetic propensities for specific traits.

Genome-wide polygenic scores (GPS) are calculated as the weighted sums of individual $i$'s SNPs:

$$GPS_{ki} = \sum_{j=i}^{m} \hat{\beta}_{kj} g_{kji}$$

Equation 2.6: Genome-wide polygenic score estimate

$GPS_{ki}$ represents the individual i's genome-wide polygenic score based on summary statistics from GWAS$_k$. $\hat{\beta}_{kj}$ is an estimate of marker $j$'s effect size for discovery trait $k$, that is, the effect of having one more copy of the reference allele at $SNP_{kj}$. $g_{kji}$ is individual $i$'s genotype at marker $j$ for discovery GWAS $k$, coded as having 0,1, or 2 copies of the reference allele at marker $k_j$. Conventionally, the $\hat{\beta}_{kj}$ for $SNP_j$ is simply the GWAS$_k$ estimate for $SNP_{jk}$. However, due to local linkage disequilibrium (LD) (i.e. correlation) between SNPs, $\hat{\beta}_{kj}$ captures any effects of the $SNP_{kj}$ and its correlates. Therefore, to correct for the multiple counting problem of effectively counting the effects of markers that are in LD with other markers multiple times, conventionally, markers are thinned down via the process of 'clumping' to a set of uncorrelated markers prior to polygenic score creation. The clumping algorithm preferentially selects the most significant markers identified by the GWAS.

Recently, several adjustments to the generic model have been proposed with the aim of improving prediction accuracy of polygenic scores. For example, LDpred attempts to avoid a reduction in predictive accuracy and loss of information caused by the conventional approach of LD-based marker pruning and applying a P-value threshold to association statistics (Vilhjalmsson et al., 2015). LDpred is a Bayesian approach that infers the posterior mean effect size of each marker by adjusting the effect size from the discovery GWAS using a prior on effect size and information on the LD between the SNPs from a reference panel to obtain a posterior estimate of the causal effect for $SNP_{jk}$ independent of the effects of other SNPs. Hence, the LDpred GPS for individual $i$ for GWAS$_k$ is the sum of $i$'s genotypes across all SNPs used in the analyses, weighted by the LDpred estimates of the genotype effects. The score represents an estimate of the genetic propensity for individual $i$ for trait $k$.

Another example is PRSice, which works on the GWAS P-value threshold for SNP inclusion into the polygenic score (Euesden et al., 2015). This method simply runs a large series of regression models and then selects the model with the 'best-fit' P-

value threshold, 'best-fit' being the score that predicts the target phenotype with the highest statistical significance. The increased multiple testing burden is addressed by an adjusted $\alpha$ level.

# 3    Phenome-wide analysis of genome-wide polygenic scores

Krapohl, E., Euesden, J., Zabaneh, D., Pingault, J.-B., Rimfeld, K., von Stumm,
S., Dale, P. S., Breen, G., O'Reilly, P. F., and Plomin, R. (2016). Phenome-wide
analysis of genome-wide polygenic scores. *Molecular Psychiatry*, 21(9):1188–1193

Supplementary material: `http://www.nature.com/mp/journal/v21/n9/suppinfo/`
`mp2015126s1.html`

# ORIGINAL ARTICLE

# Phenome-wide analysis of genome-wide polygenic scores

E Krapohl[1], J Euesden[1], D Zabaneh[1], J-B Pingault[1,2], K Rimfeld[1], S von Stumm[3], PS Dale[4], G Breen[1], PF O'Reilly[1] and R Plomin[1]

Genome-wide polygenic scores (GPS), which aggregate the effects of thousands of DNA variants from genome-wide association studies (GWAS), have the potential to make genetic predictions for individuals. We conducted a systematic investigation of associations between GPS and many behavioral traits, the behavioral phenome. For 3152 unrelated 16-year-old individuals representative of the United Kingdom, we created 13 GPS from the largest GWAS for psychiatric disorders (for example, schizophrenia, depression and dementia) and cognitive traits (for example, intelligence, educational attainment and intracranial volume). The behavioral phenome included 50 traits from the domains of psychopathology, personality, cognitive abilities and educational achievement. We examined phenome-wide profiles of associations for the entire distribution of each GPS and for the extremes of the GPS distributions. The cognitive GPS yielded stronger predictive power than the psychiatric GPS in our UK-representative sample of adolescents. For example, education GPS explained variation in adolescents' behavior problems (~0.6%) and in educational achievement (~2%) but psychiatric GPS were associated with neither. Despite the modest effect sizes of current GPS, quantile analyses illustrate the ability to stratify individuals by GPS and opportunities for research. For example, the highest and lowest septiles for the education GPS yielded a 0.5 s.d. difference in mean math grade and a 0.25 s.d. difference in mean behavior problems. We discuss the usefulness and limitations of GPS based on adult GWAS to predict genetic propensities earlier in development.

*Molecular Psychiatry* (2016) **21,** 1188–1193; doi:10.1038/mp.2015.126; published online 25 August 2015

## INTRODUCTION

One of the most striking findings emerging from genome-wide association studies (GWAS) of complex traits is the scarcity of common single nucleotide polymorphism (SNP) associations that account for more than 1% of trait variation in the population.[1,2] Although GWAS have been successful in discovering and replicating SNP associations for many traits and disorders,[3] the dearth of larger SNP associations in well-powered GWAS demonstrates that the ubiquitous heritability of complex dimensions and common disorders is caused by thousands of common DNA variants of small effect.[1,4] Because their effects are miniscule, a single common SNP is of little use for prediction. For this reason, the future of genetic prediction lies with polygenic scores that aggregate the effects of thousands of SNPs discovered by GWAS, including variants that do not achieve genome-wide significance.[5] Unlike quantitative genetic designs that estimate the net effect of DNA differences in a population—such as twin and adoption studies and SNP-based heritability[6]—polygenic scores provide individual-specific estimates of genetic propensities for specific SNPs.

Here we refer to polygenic scores as genome-wide polygenic scores (GPS) for two reasons. First, the acronym GPS excludes the term 'risk', in contrast to the previous labels, which imply that genetic influences are inevitably associated with negative outcomes. Second, the acronym GPS in its original use as global positioning system is an apt metaphor for the use of DNA differences across the genome to guide research on genetic influence.

Association statistics for dozens of large meta-analytic GWAS are now available, including GWAS for psychiatric and cognitive traits. The GPS based on these GWAS results are limited by the 'hidden heritability' ceiling and, as yet, they account for only a few percent of the variance or liability of their target trait.[2] In addition, most GWAS are based on comparisons between diagnosed cases versus controls using a liability model that assumes continuous liability throughout the population, but the extent to which these case/control results generalize to prediction of continuous traits in the population needs to be established empirically.

Multivariate quantitative genetic analyses using the twin method as well as well as SNP heritability methods have shown that genetic effects are to a substantial extent pleiotropic across complex traits in general[7] and in particular across cognitive abilities and disabilities[8,9] and across psychopathologies.[10–13] This pleiotropy suggests the usefulness of going beyond 'candidate-phenotype' analyses of a single GPS-trait pairing to consider the multivariate profile of GPS associations across many behavioral traits, the behavioral phenome.

Here, we report the first phenome-wide analysis of GPS derived from 13 published major psychiatric, cognitive and biometric GWAS. We applied effect size and significance estimates from GWAS summary statistics to create GPS from raw genotype data for individuals in our target sample. The phenome included 50 traits from the domains of psychopathology, personality, cognitive abilities, and educational achievement, assessed in a representative sample of over 3000 16-year-old individuals in the United Kingdom.

[1]MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, UK; [2]Division of Psychology and Language Sciences, University College London, London, UK; [3]Department of Psychology, Goldsmiths University of London, New Cross, London, UK and [4]Department of Speech and Hearing Sciences, University of New Mexico, Albuquerque, NM, USA. Correspondence: Professor R Plomin, MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology and Neuroscience, King's College London, DeCrespigny Park, Denmark Hill, London SE5 8AF, UK.
E-mail: robert.plomin@kcl.ac.uk

The main focus of this paper is to explore the profile of GPS associations across the behavioral phenome for the entire distribution of each GPS and for the extremes of the GPS distributions. One use of polygenic scores is to predict genetic propensities early in development in order to facilitate interventions that promote potential and prevent problems. As a step in this direction, the present sample consists of adolescents as they finish compulsory schooling at age 16. We test whether GPS, based on current GWAS, predict phenotypic variation in the adolescent population, and we discuss the usefulness and limitations of GPS based on adult GWAS to predict genetic propensities earlier in development.

## MATERIALS AND METHODS
We used genome-wide genotype and phenome-wide behavioral data from 3152 unrelated adolescents drawn from the UK-representative Twins Early Development Study[14–16] (Supplementary Table 1). We processed the 3152 genotypes using standard quality control procedures followed by imputation of SNPs using the 1000 Genomes Project reference panel[17] (Supplementary Methods 1). After quality control, we included around 4.3 million variants into the polygenic score analyses (Supplementary Methods 1). Association analyses were conducted using imputed markers and principal components to control for population stratification. Individuals were assessed on a wide range of phenotypes at the age of 16. The present analyses included 50 traits from the domains of psychopathology, personality, cognitive abilities and educational achievement (Supplementary Methods 2). All measures were age- and sex-regressed and the z-scores were used in the analyses.

We created 13 GPS for each of the over 3000 individuals in our sample using summary statistics from 13 published GWAS[18–28] (Supplementary Table 2). Here we present results using a $P = 0.30$ threshold for including SNPs from the published GWAS (Figure 1 and Supplementary Table 3); results for GPS based on the $P$-value thresholds of 0.10 and 0.05 are included in the Supplementary material (Supplementary Figures 1a and b and Supplementary Table 3). The selection of the relatively lenient $P = 0.30$ threshold was based on the evidence that many associated markers lie within the ensemble of individually non-significant SNPs, with power of the

GPS increasing with number of SNP included.[5] We also report results (Supplementary Figure 2 and Supplementary Table 4) from a high-resolution polygenic scoring approach, implemented in the software PRSice (London, UK), that identifies the most predictive GPS for each phenotype.[29]

We describe two types of main results: (i) associations between GPS and the behavioral phenome for the entire sample, which demonstrate the usefulness of cross-trait prediction, and (ii) quantile analyses showing the association between selected GPS and behavior by septile, which illustrates the ability to stratify individuals by GPS and the potential of polygenic score for phenotype prediction.

To inform these analyses, we demonstrate that GPS are normally distributed and discuss the implications for considering both ends—resilience as well as risk—of GPS distributions. We also examine three types of correlations: (i) genetic correlations between the GWAS summary statistics (ii) correlations between the GPS, and (iii) phenotypic correlations between the target phenotypes. These correlations support the usefulness of a phenome-wide analysis of GPS.

## RESULTS

### GPS are normally distributed
The quantitative genetic model assumes that many genetic variants of small effect drive the heritability of complex traits and common disorders,[30] even though each marker is inherited in the discrete manner hypothesized by Mendel.[31] Therefore, the central limit theorem implies that the distribution of polygenic scores in the population will approach normality. Specifically, the normal distribution is to be expected whenever trait variation is polygenic and produced by the addition of a large number of small effects.

Nonetheless, the normality of GPS (Supplementary Figures 3a and c) merits emphasis because it illustrates that common disorders can be considered as extremes of the common polygenic liability spectrum, which has far-reaching implications for diagnosis, treatment and prevention.[32] It also implies that GPS can be operationalized in terms of 'resilience' as well as 'risk' predictors. There is untapped research potential for operationalizing the



**Figure 1.** Correlations between 13 genome-wide polygenic scores and 50 traits from the behavioral phenome. These results are based on GPS constructed using a GWAS $P$-value threshold ($P_T$) = 0.30; results for $P_T$ = 0.10 and 0.05 (Supplementary Figures 1a and b and Supplementary Table 3). $P$-values that pass Nyholt–Sidak correction (see Supplementary Methods 1) are indicated with two asterisks, whereas those reaching nominal significance (thus suggestive evidence) are shown with a single asterisk.

negative tail of GPS for disorders as 'resilience' and the negative end of cognitive or education GPS as 'risk' factors. This 'other end' of the normal distribution of GPS is uncharted territory. From an evolutionary perspective, averageness might be an adaptive trade-off against the mishmash of costs and benefits of more extreme GPS, especially given the fluctuating nature of selection.[32]

### Intercorrelations between GWAS, GPS and phenotypic traits

As depicted in Supplementary Figure 4 the phenotypic correlations between the target phenotypes in our sample of adolescents show substantial intercorrelations, with a 'cognitive' and a 'psychopathology' cluster.

We estimate the genetic correlation between the discovery GWAS using a new technique based on LD score regression,[33,34] which uses only GWAS summary statistics and linkage disequilibrium information to decompose true polygenic variance/ covariance from confounding (see Supplementary Methods for details). Supplementary Figure 5 depicts the genetic correlations between the 13 GWAS, which provide evidence for significant and substantial pleiotropy. In addition to the genetic correlations reported previously,[34] we add correlations for the summary statistics of the child IQ GWAS,[19] adult IQ[35] and intracranial volume.[20] The observed genetic correlations replicate and extend previous research. We confirm genetic overlap between the major psychoses[13,25,34,36,37] and between cognitive phenotypes including intracranial volume,[18,20,38–40] respectively. We further find correlations between these two clusters—for example, strong negative associations between the cognitive phenotypes and Alzheimer's and positive associations between educational attainment and autism spectrum disorder as well as bipolar disorder.

We also examined correlations between the GPS created for our sample (Supplementary Figures 6a and c and Supplementary Table 3). We find similar correlation patterns but weaker overall correlations.

These genetic correlations provide evidence that polygenic effects are to a substantial degree pleiotropic across traits. Together with finding substantial correlations between the target phenotypes, this multivariate genetic architecture suggests the usefulness of a phenome-wide approach to investigate the links between GPS and behavior, which is the focus of the next and final section of results.

### GPS correlate with the behavioral phenome

Figure 1 summarizes correlations between the 50 traits of the behavioral phenome and the 13 GPS for $P_T = 0.30$. Correlation coefficients, s.e., $P$-value thresholds ($P_T$), and number of SNPs included are shown in Supplementary Table 3 for the fixed $P_T$ (0.30; 0.10; 0.05). Very similar patterns of association emerged from both the conventional fixed $P_T$ analyses and the high-resolution analyses that estimate the $P_T$ flexibly for the 'best-fit' GPS (Supplementary Figure 2 and Supplementary Table 4). Both methods yielded statistically significant phenomic associations only for the GPS for College and Child IQ.

*College GPS.* College GPS, which was based on the binary measure of attending college or not, showed the strongest phenomic profile at age 16, which might reflect the fact that its meta-analytic GWAS sample size was one of the largest ($N = 120 000$; Rietveld *et al.*[18]). College GPS correlated significantly with academic performance at age 16: General Certificate of Secondary Examination (GCSE) English ($r = 0.15$), GCSE mathematics ($r = 0.15$, s.e. 0.02) and GCSE science ($r = 0.14$, s.e. 0.02).[39] College GPS also correlated significantly with general cognitive ability ('g') ($r = 0.14$, s.e. 0.03) as well as its subscales Ravens Matrices ($r = 0.12$, s.e. 0.03) and with Mill Hill Vocabulary ($r = 0.09$, s.e. 0.03), which confirms a similar finding for adults.[40] College GPS also correlated positively with PISA math interest ($r = 0.10$, s.e. 0.03) and math

self-efficacy ($r = 0.12$, s.e. 0.03). Negative associations for College GPS emerged for SDQ total behavior problems ($r = -0.07$, s.e. 0.02) and SDQ Conduct ($r = -0.08$, s.e. 0.02).

*Child IQ GPS.* The GPS for Child IQ yielded a similar but diluted phenomic profile as compared with College GPS. Child IQ GPS correlated significantly with GCSE English ($r = 0.09$, s.e. 0.02), GCSE Math ($r = 0.10$, s.e. 0.02) and GCSE Science ($r = 0.09$, s.e. 0.02).

*Psychiatric GPS.* In contrast, the five psychiatric GPS yielded no significant correlations that passed multiple comparisons corrections across the behavioral phenome. Nominally significant associations included a positive correlation between Alzheimer's GPS and Conner's Impulsivity; and positive associations between Autism Spectrum Disorder GPS and Autism Quotient: Attention Switching. Autism Spectrum Disorder GPS yielded nominally significant negative associations with Chaos at home, Attachment and Height. Schizophrenia GPS correlated positively with GCSE English and negatively with Autism Quotient: Attention to Detail. Bipolar disorder GPS correlated negatively with Autism Quotient: Attention to Detail.

One likely explanation for the lower phenomic profile of psychiatric GPS compared with that of College GPS is the difference in sample sizes for the discovery samples. However, Child IQ GPS yielded significant associations despite the relatively smaller sample size of the GWAS ($N = 9616$). This might point to the importance of developmental proximity or similarity of the phenotypes in discovery and target sample. It also emphasizes that predictive power is not only a function of sample size of the discovery sample.[5] Phenotypic similarity between the traits in the discovery sample and the target sample is a proxy for the magnitude of genetic covariance between the traits.

The underlying premise of GWAS is that the polygenic architecture of complex traits and common disorders requires a genome-wide approach despite the multiple testing burden. Similarly, based on strong evidence for the ubiquitous pleiotropy of complex traits,[7,9–13,34] the advantage of the phenome-wide approach outweighs the resulting multiple testing burden. Specifically, while testing a large number of highly unlikely hypotheses with little or no prior support should be avoided, in this case we have collated a well-defined set of psychological and behavioral traits for which there is good reason to suspect causal associations with the available discovery GWAS phenotypes. In this way, the only 'multiple testing problem' relates to setting an appropriate significance threshold given the number and correlation of tests performed (see Supplementary Methods for multiple testing correction method used).

Therefore, the absence of phenome-wide significant associations (that is, after correcting for multiple testing across the 50 traits and 13 GPS) for all psychiatric GPS does not imply the absence of polygenic effects. However, the scarcity of nominally significant associations between the psychiatric GPS and the 50 traits suggests that the genetic covariance between psychiatric adult case/control samples and our adolescent population sample might be relatively small. For instance, under certain assumptions about polygenic architecture (for example, $\leqslant 5\%$ of tested SNPs associated with schizophrenia in the discovery GWAS), we had $\geqslant 80\%$ power with $\alpha = 0.05$ to detect associations between the Schizophrenia GPS and a phenotype given $\geqslant 0.06$ genetic covariance between schizophrenia and the target trait, with $\geqslant 0.5\%$ of phenotypic variation in the target trait explained by schizophrenia[5,41,42] (see Supplementary Methods for more detail).

One possible reason for the lower observed phenomic profile of the psychiatric GPS might be that the current sample is UK representative and therefore not enriched for psychiatric symptoms. The psychiatric GPS were based on case–control comparisons, often with extreme cases. This emphasizes the limitations of using GPS for

**Figure 2.** (**a**) Mean for height at age 16 by adult Height genome-wide polygenic score (GPS) septile. The threshold for selecting trait-associated alleles was $P_T < 0.30$. The GPS were converted to quantiles (1 = lowest, 7 = highest GPS). Mean phenotypic values and 95% confidence intervals (CIs) for the quantile groups (bars) were estimated using general linear regression with ancestrally informative principal components, sex and age of measurement as covariates. (**b**) Mean for children's mathematics educational achievement at age 16 (compulsory subject on the General Certificate of Secondary Examination (GCSE), see Materials and Methods for details) by College GPS septile. The threshold for selecting trait-associated alleles was $P_T < 0.30$. The GPS were converted to quantiles (1 = lowest, 7 = highest GPS). Mean phenotypic values and 95% CI for the quantile groups (bars) were estimated using general linear regression with ancestrally informative principal components, sex and age of measurement as covariates. (**c**) Mean for total parent-reported behavior problems at age 16 by adult College GPS septile. The threshold for selecting trait-associated alleles was $P_T < 0.30$ (the best-fit GPS as estimated by PRSice software, see Materials and Methods). The GPS were converted to quantiles (1 = lowest, 7 = highest GPS). Mean phenotypic values and 95% CI for the quantile groups (bars) were estimated using general linear regression with ancestrally informative principal components, sex and age of measurement as covariates.

the prediction of trait variation in the general population from GWAS based on selected samples. Importantly, the GPS College did predict children's behavior problems in our UK-representative sample, whereas the psychiatric GPS did not. This points to the usefulness of cross-trait prediction in general and the value of cognitive GWAS/GPS as prediction instruments for psychiatric symptoms in the population.

*Other GPS.* Adult body mass index (BMI) GPS correlate positively with the measure of BMI at age 16 ($r = 0.18$, s.e. 0.03); and adult Height GPS correlate with height at age 16 ($r = 0.33$, s.e. 0.03). There was suggestive evidence for a negative association between Ever smoked GPS and conscientiousness ($r = -0.06$, s.e. 0.03) and a positive association with BMI ($r = 0.09$, s.e. 0.03).

Quantile analyses
To illustrate the ability to stratify individuals by GPS and the potential of polygenic score for phenotype prediction, we grouped individuals into GPS septiles and estimated the mean phenotypic value for each quantile. We provide three examples:
Figure 2a shows that mean standardized height increased with more adult height-associated alleles in our UK-representative sample of children aged 16, with the largest difference between the lowest and highest septile (Hedges *g*: − 1.01 with 95% confidence interval (CI): − 1.26 to − 0.77; difference in means: 0.97 s.d., with *P*-value < 0.01). Figure 2b shows that mean math grade on the standardized UK-national examinations at age 16 increased with more College-associated alleles (Figure 2b), with the largest difference between the lowest and highest septile (Hedges *g*: − 0.52 with 95% CI: − 0.67 to − 0.37; difference in means: 0.49 s.d., with *P*-value < 0.01).

Figure 2c illustrates the utility of the phenome-wide approach for cross-trait prediction: the mean for total parent-reported behavior problems at age 16 decreased slightly but significantly with higher College GPS, with a maximum effect size between the lowest and highest quantile (Hedges *g*: 0.20 with 95% CI: 0.04–0.34; difference in means: 0.19, with *P*-value 0.01).
These results (Figure 2) illustrate the ability to stratify individuals by GPS, which suggests opportunities for research, for example, selecting high and low GPS extreme individuals for intensive research such as neuroimaging that is unable to test large representative samples. However, we emphasize that the current predictive power and accuracy of GPS do not allow for their use as predictive tests.

**DISCUSSION**
These results highlight the usefulness of a phenome-wide approach to examine behavioral profiles of associations with GPS even though current GPS account for only a few percent of variance or liability of their target trait. An interesting finding is that phenome-wide associations for cognitive GPS are stronger than for psychiatric GPS in our UK-representative sample of adolescents. For example, we found that GPS College, but none of the psychiatric GPS, predicted adolescent behavior problems, which demonstrates the usefulness of cross-trait predictions and the multivariate phenome-wide approach in general. However, this finding could be explained by differences in sample sizes, sampling methods (population versus case/control), and genetic architecture (for example, extent of covariance between discovery and target trait).

Finding significant associations for the Child IQ GPS, which is based on a small discovery sample, is a reminder that predictive power of GPS is not merely a function of sample size but also of the developmental proximity of the GWAS sample and the target GPS sample. As explained in the Introduction, we were interested in the extent to which GWAS in adult samples yield GPS that can predict genetic propensities—strengths as well as weaknesses—earlier in development, in this case in adolescence. However, GPS College is a trait assessed closer in age to the adolescents in our sample. In contrast, the psychiatric GPS were derived from GWAS studies of adults.

A larger issue is that extant GPS account for only a few percent of the phenomic variance in the target trait. However, we illustrate the research potential of polygenic stratification by quantile. Power and accuracy of GPS will improve as GWAS sample sizes increase. GPS that narrow the 'hidden heritability' gap is what is needed most for phenome-wide analyses—and for all research harvesting the fruits of GWAS.

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## AUTHOR CONTRIBUTIONS

RP directs and received funding for the Twins Early Development Study (TEDS). RP and EK conceived of the present study. EK analyzed and interpreted the data. RP and PFO supervised the project and interpreted the data. RP and EK wrote the manuscript with help from PFO, DZ, J-BP, KR, SvS, PSD and GB. JE provided an early version of the PRSice software (http://PRSice.info/).

## REFERENCES

1 Gratten J, Wray NR, Keller MC, Visscher PM. Large-scale genomics unveils the genetic architecture of psychiatric disorders. *Nat Neurosci* 2014; **17**: 782–790.

2 Wray NR, Lee SH, Mehta D, Vinkhuyzen AAE, Dudbridge F, Middeldorp CM. Research review: polygenic methods and their application to psychiatric traits. *J Child Psychol Psychiatry* 2014; **55**: 1068–1087.

3 Visscher PM, Brown MA, McCarthy MI, Yang J. Five years of GWAS discovery. *Am J Hum Genet* 2012; **90**: 7–24.

4 Robinson MR, Wray NR, Visscher PM. Explaining additional genetic variation in complex traits. *Trends Genet* 2014; **30**: 124–132.

5 Dudbridge F. Power and predictive accuracy of polygenic risk scores. *PLoS Genet* 2013; **9**: e1003348.

6 Yang JA, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet* 2011; **88**: 76–82.

7 Sivakumaran S, Agakov F, Theodoratou E, Prendergast JG, Zgaga L, Manolio T *et al*. Abundant pleiotropy in human complex diseases and traits. *Am J Hum Genet* 2011; **89**: 607–618.

8 Plomin R, Kovas Y. Generalist genes and learning disabilities. *Psychol Bull* 2005; **131**: 592–617.

9 Trzaskowski M, Davis OSP, DeFries JC, Yang J, Visscher PM, Plomin R. DNA evidence for strong genome-wide pleiotropy of cognitive and learning abilities. *Behav Genet* 2013; **43**: 267–273.

10 Kendler KS, Aggen SH, Knudsen GP, Røysamb E, Neale MC, Reichborn-Kjennerud T. The structure of genetic and environmental risk factors for syndromal and subsyndromal common DSM-IV Axis I and All Axis II disorders. *Am J Psychiatry* 2011; **168**: 29–39.

11 Lichtenstein P, Carlström E, Råstam M, Gillberg C, Anckarsäter H. The genetics of autism spectrum disorders and related neuropsychiatric disorders in childhood. *Am J Psychiatry* 2010; **167**: 1357–1363.

12 Lichtenstein P, Yip BH, Björk C, Pawitan Y, Cannon TD, Sullivan PF *et al*. Common genetic determinants of schizophrenia and bipolar disorder in Swedish families: a population-based study. *Lancet* 2009; **373**: 234–239.

13 Cross-Disorder Group of the Psychiatric Genomics Consortium. Genetic relationship between five psychiatric disorders estimated from genome-wide SNPs. *Nat Genet* 2013; **45**: 984–994.

14 Haworth CMA, Davis OSP, Plomin R. Twins Early Development Study (TEDS): a genetically sensitive investigation of cognitive and behavioral development from childhood to young adulthood. *Twin Res Hum Genet* 2013; **16**: 117–125.

15 Oliver BR, Plomin R. Twins' Early Development Study (TEDS): a multivariate, longitudinal genetic investigation of language, cognition and behavior problems from childhood through adolescence. *Twin Res Hum Genet* 2007; **10**: 96–105.

16 Kovas Y, Haworth CMA, Dale PS, Plomin R. The genetic and environmental origins of learning abilities and disabilities in the early school years. *Monogr Soc Res Child Dev* 2007; **72**, vii 1–144.

17 Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE *et al*. An integrated map of genetic variation from 1,092 human genomes. *Nature* 2012; **491**: 56–65.

18 Rietveld CA, Medland SE, Derringer J, Yang J, Esko T, Martin NW *et al*. GWAS of 126,559 individuals identifies genetic variants associated with educational attainment. *Science* 2013; **340**: 1467–1471.

19 Benyamin B, Pourcain BS, Davis OS, Davies G, Hansell NK, Brion M-J *et al*. Childhood intelligence is heritable, highly polygenic and associated with FNBP1L. *Mol Psychiatry* 2014; **19**: 253–258.

20 Hibar DP, Stein JL, Renteria ME, Arias-Vasquez A, Desrivières S, Jahanshad N *et al*. Common genetic variants influence human subcortical brain structures. *Nature* 2015; **520**: 224–229.

21 Lambert JC, Ibrahim-Verbaas CA, Harold D, Naj AC, Sims R, Bellenguez C *et al*. Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. *Nat Genet* 2013; **45**: 1452–1458.

22 Psychiatric GWAS Consortium Bipolar Disorder Working Group. Large-scale genome-wide association analysis of bipolar disorder identifies a new susceptibility locus near ODZ4. *Nat Genet* 2011; **43**: 977–983.

23 Major Depressive Disorder Working Group of the Psychiatric GWAS Consortium, Ripke S, Wray NR, Lewis CM, Hamilton SP, Weissman MM *et al*. A mega-analysis of genome-wide association studies for major depressive disorder. *Mol Psychiatry* 2013; **18**: 497–511.

24 Schizophrenia Working Group of the Psychiatric Genomics Consortium. Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 2014; **511**: 421–427.

25 Cross-Disorder Group of the Psychiatric Genomics Consortium. Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *Lancet* 2013; **381**: 1371–1379.

26 Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S *et al*. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat Genet* 2014; **46**: 1173–1186.

27 Tobacco and Genetics Consortium. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nat Genet* 2010; **42**: 441–447.

28 Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR *et al*. Genetic studies of body mass index yield new insights for obesity biology. *Nature* 2015; **518**: 197–206.

29 Euesden J, Lewis CM, O'Reilly PF. PRSice: polygenic risk score software. *Bioinformatics* 2015; **31**: 1466–1468.

30 Fisher RA. The correlation between relatives on the supposition of Mendelian inheritance. *Transactions of the Royal Society of Edinburgh* 1918; **52**: 399–433.

31 Mendel G. Versuche über Pflanzenhybriden. *Verhandlungen Naturforschenden Vereines Brunn* 1866; **4**: 3–47.

32 Plomin R, Haworth CMA, Davis OSP. Common disorders are quantitative traits. *Nat Rev Genet* 2009; **10**: 872–878.

33 Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al*. LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat Genet* 2015; **47**: 291–295.

34 Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, ReproGen Consortium *et al*. An atlas of genetic correlations across human diseases and traits. *bioRxiv* 2015; 014498.

35 Davies G, Armstrong N, Bis JC, Bressler J, Chouraki V, Giddaluru S *et al*. Genetic contributions to variation in general cognitive function: a meta-analysis of genome-wide association studies in the CHARGE consortium (N = 53 949). *Mol Psychiatry* 2015; **20**: 183–192.

36 Kavanagh DH, Tansey KE, O'Donovan MC, Owen MJ. Schizophrenia genetics: emerging themes for a complex disorder. *Mol Psychiatry* 2015; **20**: 72–76.

37 Maier R, Moser G, Chen G-B, Ripke S, Cross-Disorder Working Group of the Psychiatric Genomics Consortium, Coryell W *et al.* Joint analysis of psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major depressive disorder. *Am J Hum Genet* 2015; **96**: 283–294.

38 Trzaskowski M, Harlaar N, Arden R, Krapohl E, Rimfeld K, McMillan A *et al.* Genetic influence on family socioeconomic status and children's intelligence. *Intelligence* 2014; **42**: 83–88.

39 Krapohl E, Plomin R. Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs. *Mol Psychiatry* advance online publication, 10 March 2015; doi:10.1038/mp.2015.2.

40 Rietveld CA, Esko T, Davies G, Pers TH, Turley P, Benyamin B *et al.* Common genetic variants associated with cognitive performance identified using the proxy-phenotype method. *Proc Natl Acad Sci USA* 2014; **111**: 13790–13794.

41 Wray NR, Goddard ME, Visscher PM. Prediction of individual genetic risk to disease from genome-wide association studies. *Genome Res* 2007; **17**: 1520–1528.

42 Chatterjee N, Wheeler B, Sampson J, Hartge P, Chanock SJ, Park J-H. Projecting the performance of risk prediction based on polygenic analyses of genome-wide association studies. *Nat Genet* 2013; **45**: 400–405.

Supplementary Information accompanies the paper on the Molecular Psychiatry website (http://www.nature.com/mp)

# 4 Multi-polygenic score approach to trait prediction

Krapohl, E., Patel, H., Newhouse, S., Curtis, C., von Stumm, S., Dale, P., Zabaneh, D., Breen, G., O'Reilly, P., and Plomin, R. (2017b). Multi-polygenic score approach to trait prediction. *Molecular psychiatry*

Supplementary material: `https://www.nature.com/mp/journal/vaop/ncurrent/suppinfo/mp2017163s1.html?url=/mp/journal/vaop/ncurrent/full/mp2017163a.html`

## ORIGINAL ARTICLE

# Multi-polygenic score approach to trait prediction

E Krapohl[1], H Patel[2,3], S Newhouse[2,3,4], CJ Curtis[1,2], S von Stumm[5], PS Dale[6], D Zabaneh[1], G Breen[1,2], PF O'Reilly[1] and R Plomin[1]

A primary goal of polygenic scores, which aggregate the effects of thousands of trait-associated DNA variants discovered in genome-wide association studies (GWASs), is to estimate individual-specific genetic propensities and predict outcomes. This is typically achieved using a single polygenic score, but here we use a multi-polygenic score (MPS) approach to increase predictive power by exploiting the joint power of multiple discovery GWASs, without assumptions about the relationships among predictors. We used summary statistics of 81 well-powered GWASs of cognitive, medical and anthropometric traits to predict three core developmental outcomes in our independent target sample: educational achievement, body mass index (BMI) and general cognitive ability. We used regularized regression with repeated cross-validation to select from and estimate contributions of 81 polygenic scores in a UK representative sample of 6710 unrelated adolescents. The MPS approach predicted 10.9% variance in educational achievement, 4.8% in general cognitive ability and 5.4% in BMI in an independent test set, predicting 1.1%, 1.1%, and 1.6% more variance than the best single-score predictions. As other relevant GWA analyses are reported, they can be incorporated in MPS models to maximize phenotype prediction. The MPS approach should be useful in research with modest sample sizes to investigate developmental, multivariate and gene–environment interplay issues and, eventually, in clinical settings to predict and prevent problems using personalized interventions.

*Molecular Psychiatry* advance online publication, 8 August 2017; doi:10.1038/mp.2017.163

## INTRODUCTION

Genome-wide association studies (GWASs) have been successful in identifying thousands of associations for hundreds of complex traits and common disorders.[1] One use of GWAS results is to understand biological pathways between genotypes and phenotypes. Another use, the focus of the present research, is to estimate genetic propensities of individuals to predict individuals' future problems and potential and, eventually, to develop personalized interventions that meet individual medical, psychiatric and educational needs. Both goals have been hindered by the ubiquitous GWA finding that the largest effect sizes are extremely small.[2] For example, the largest population effect sizes found for common variants in height or body mass index (BMI) account for only ~1% of the variance.[3,4] We know empirically that the vast majority of common genetic variants for most traits have a markedly lower effect than 1%.[2]

The highly polygenic nature of complex traits and common disorders poses an immense challenge for understanding the biological mechanisms linking single variants with phenotypes. However, when the priority is phenotypic prediction, polygenic scores can be used to aggregate the effects of many DNA variants in order to investigate their joint predictive power.[5,6] Rather than just using single-nucleotide polymorphisms (SNPs) that reach genome-wide significance, a recent development is to aggregate a much larger number of SNPs, weighted by their GWA effect size estimate, as long as together they increase the prediction in an independent sample, even if some SNPs have no real effect.[7] For example, for height, a polygenic score that aggregates the effects of ~2000 SNPs accounts for 21% of the variance of height in independent samples.[3]

The other defining characteristic of complex traits and common disorders is the abundance of genetic correlations between them. There is consistent evidence for genetic correlations between psychiatric disorders, between anthropometric traits and between educational and cognitive traits, as well as for genetic correlations across these categories.[8–11]

Genetic correlation can arise from pleiotropy, the phenomenon of multiple traits being associated with the same gene or genetic variant.[8] Genetic correlation can also reflect shared biological pathways or more indirect linkage.[12] Regardless of its cause, genetic correlation between different traits means that a polygenic score based on one trait can predict a different outcome trait, with predictive accuracy a function of the shared genetic signal between them. Therefore, when the aim is prediction, genetic correlation can be exploited for trait prediction while remaining agnostic to the underlying mechanisms.

A primary goal of polygenic scores, which aggregate the effects of thousands of trait-associated genetic variants discovered in GWAS, is to estimate individual-specific genetic propensities. This is typically achieved using a single polygenic score, but here we use an approach to increase predictive power by exploiting the joint power of multiple discovery GWASs. We use a multi-polygenic score (MPS) approach that exploits genetic correlations between the outcome trait and a multitude of traits by using the joint predictive power of multiple polygenic scores in one regression model.

[1]MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, UK; [2]Department of Biostatistics and Health Informatics, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, UK; [3]NIHR Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and King's College London, London, UK; [4]Farr Institute of Health Informatics Research, UCL Institute of Health Informatics, University College London, London, UK; [5]Department of Psychology, Goldsmiths University of London, New Cross, London, UK and [6]Department of Speech and Hearing Sciences, University of New Mexico, Albuquerque, NM, USA. Correspondence: E Krapohl, MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology and Neuroscience, King's College London, 16 DeCrespigny Park, London SE5 8AF, UK.
E-mail: eva.krapohl@kcl.ac.uk
Received 19 October 2016; revised 12 May 2017; accepted 20 June 2017

We selected GWASs from a centralized repository of summary statistics—based on their statistical power and regardless of prior evidence for association with the outcomes—to predict three core developmental outcomes in our independent target sample: educational achievement, BMI, and general cognitive ability. Using repeated cross-validation, we trained and validated the prediction models using elastic net regularized regression, a multiple regression model suited to deal with a large number of correlated predictors while preventing overfitting.[13] We subsequently tested how well these models predict outcomes in an independent test set.

Here, we employ a MPS approach that uses publicly available GWAS summary statistics to estimate individual-level genetic propensities and predict developmental outcomes in an independent target sample. This stands in contrast to multi-trait approaches that rely on access to individual-level data in the discovery data sets because they make use of a method from animal breeding in which the total genetic effect ('breeding value') of each individual in a discovery data set is estimated from the best linear unbiased predictor in a multi-trait random-effects model that can be used for individual-level prediction in the validation data sets. These multi-trait methods are not applicable to GWAS summary statistics when genotype data are unavailable because of privacy or logistical constraints that are frequently the case.

The declared aim of the current MPS approach is to maximize prediction of developmental outcomes, rather than investigating their etiology. This stands in contrast to multi-trait meta-analytic approaches of GWAS summary statistics that relies on substantial and consistent correlations between discovery GWASs and whose main aim is variant discovery.[14–17] The current MPS approach allows for, but does not require, correlation among polygenic predictors.

## MATERIALS AND METHODS

### Sample

The target sample comprised genome-wide SNP and phenotypic data from 6710 unrelated adolescents drawn from the UK representative Twins Early Development Study (TEDS). TEDS is a multivariate longitudinal study that recruited over 11 000 twin pairs born in England and Wales in 1994, 1995 and 1996. Both the overall TEDS sample and the genotyped subsample have been shown to be representative of the UK population.[18–20] The project received approval from the Institute of Psychiatry ethics committee (05/Q0706/228) and parental consent was obtained before data collection. We processed the genotypes for the 6710 individuals using stringent quality control procedures followed by imputation of SNPs using the Haplotype Reference Consortium reference panel[21] (Supplementary Methods S1).

### Predictors

*Discovery data sets: GWAS summary statistics.* We selected GWAS summary statistics from *LD hub*, a centralized repository for summary statistics[22] based on their statistical power—regardless of prior evidence for association with our outcome traits. Specifically, we included 81 GWAS summary statistics that were either publically downloadable or obtained via correspondence and had a linkage disequilibrium (LD) score[23] heritability z-score > 5, indexing good statistical power (which is a function of variance explained and sample size). Supplementary Table S1 provides details of all GWAS summary statistics included in our analyses.

The published version of the child IQ GWAS included the present target sample of TEDS. Therefore, to avoid bias, the present analyses used summary statistics from a rerun of the GWAS meta-analysis excluding TEDS.

*Polygenic scores.* We created 81 genome-wide polygenic scores for each of the 6710 individuals in the TEDS sample using summary statistics from the GWAS described above (Supplementary Table S1). After quality control (Supplementary Methods S1), the study data included 7 581 516 genotyped or well-imputed (info > 0.70) SNPs. These were quality controlled

and coordinated with each of the summary statistics, respectively, by excluding markers due to nucleotide inconsistencies or low minor allele frequency ( < 1%). Number of markers before and after quality control and coordination with the study data are listed in Supplementary Table S1.

We constructed polygenic scores as the weighted sums of each individual's trait-associated alleles across all SNPs. We used LDpred[24] to construct the scores. LDpred uses a prior on the markers' effect sizes and adjusts summary statistics for LD between markers. Scores were standardized and adjusted for 30 principal components. More details on the construction of the polygenic scores are provided in Supplementary Methods S2.

### Outcomes

To illustrate the MPS approach, we selected three key developmental outcomes:

Educational achievement operationalized as the mean grade of the three compulsory subjects (Mathematics, English and Science) attained on the standardized United Kingdom General Certificate of Secondary Education (GCSE), taken by almost all ( > 99%) pupils at the end of compulsory education at age 16 years.

General cognitive ability at age 12 years assessed by two verbal and two nonverbal cognitive standardized tests.

BMI at age 9 years that was age and sex adjusted using external reference data.

Supplementary Methods S3 and Figure S1 contain detailed descriptions of the three measures.

### Models

*Single-polygenic score models.* To estimate the separate prediction of each predictor, we fit a series of simple linear regression models for each of the 81 polygenic scores and each of the 3 outcomes. For each GWAS-outcome combination, three models were run using polygenic scores created with Gaussian mixture weights of 1, 0.1 and 0.01, respectively. The model that explained the most variance in the outcome (that is, largest cross-validated $R^2$ in training data) was then entered into the multi-score model. These simple linear regression models were fit and validated in repeated 10-fold cross-validation (see section below for details) using the *lm* function implemented within the *caret* R package.[25] Based on consistent evidence for extensive genetic correlations across complex traits and disorders, rather than summing up, the predictions of the single-score models were expected to substantially overlap.

*MPS models.* We used the MPS model to estimate the joint prediction of the 81 polygenic scores as well as the ranking of predictors by the magnitude of their contribution to predicting the outcome.

Conventional multiple linear regression models in the presence of a large number of predictors are subject to overfitting, and stepwise regression suffers from upward-biased coefficients and $R^2$ (see, for example, Tibshirani[26]). We used elastic net regularized regression[13] to predict outcomes and by selecting predictors and estimating their contribution to the prediction. Regularized regression models are general linear models that employ strict penalties to prevent overfitting. Elastic net allows for estimating the joint predictive ability of a large number of variables while preventing overfitting. Elastic net uses a linear combination of two regularization techniques, L2 regularization (used in ridge regression) and L1 regularization (used in LASSO (least absolute shrinkage and selection operator)) by simultaneously implementing variable selection (that is, dropping/retaining variables) and continuous shrinkage (that is, penalizing coefficients for overfitting); and it efficiently deals with multicollinearity by selecting or dropping groups of correlated variables.[13,27]

Elastic net overcomes the limitation of LASSO that tends to select one variable from a group of correlated predictors and to ignore the others. In situations where predictors are non-independent or correlated (for example, sharing genetic signal or discovery cohorts) the elastic net has the advantage of including automatically all the highly correlated variables in the group (*grouping effect*).[13,27,28]

Final model coefficients are analogous to a conventional multiple linear regression output that allows for a ranking of predictors by the magnitude of their contribution to predicting the outcome. Overall variance explained by the model is indexed by the coefficient of determination, $R^2$.

We used *glmnet* R package[15–17] implemented within *caret* R package[25] to conduct a series of linear elastic net regularized regressions and select polygenic predictors leading to an optimized final model for each

outcome. Elastic net regularized regression employs two hyperparameters, alpha and lambda.[13] As recommended to achieve optimized balance between variance explained and minimum bias, we fit models to tune over both alpha and lambda parameter values in repeated 10-fold cross-validation.[29]

### Model training and testing

Generally, a predictive model is considered powerful when the model is capable of predicting outcomes in 'unseen' data with high accuracy. The performance of a model can therefore be evaluated by testing how well it predicts phenotypes of individuals whose data were not included in the construction of the prediction model.

Each model described in the preceding section was trained and tested using the following three-step strategy:

*Data splitting*. We randomly split the data set into a separate training set and test set (60% *train*, 40% *test*).

*Model training*. We used repeated cross-validation on the training set to train and optimize the model via validation.

*Model testing and comparison*. We applied the final model to the independent test set to obtain an unbiased estimate of model performance.

*Model training*.   The training set was used to train and validate the model, this included hyperparameter tuning for the elastic net models. In order to optimize the balance between variance explained and minimum bias, we tested each model in 10-fold cross-validation with resampling.[29] We split the training data randomly into 10 equal-sized subsets, using 9 subsets to train the model and the remaining subset as validation. The cross-validation process was repeated 10 times, with each of the 10 subsamples used once as the validation data.

Although cross-validation has been shown to produce nearly unbiased estimates of accuracy, variability of these estimates can be reduced by bootstrap methods, wherein available data are repeatedly sampled with replacement in order to mimic the drawing of future random sampling.[30,31] Therefore, to minimize variation across validation data sets, we repeated the 10-fold cross-validation 100 times with random data set partitions.[32]

The optimized or 'final' model is chosen based on the largest performance value (or smallest mean squared error). Predictors retained within the model and standardized coefficients index whether, and to what extent, they contribute to predicting the outcome. Model performance for the repeated cross-validation in the training set was summarized as mean-cv-$R^2_{train}$ from the resampling distribution.

*Model testing and comparison*. To obtain unbiased estimates of model performance, we used the parameters from the final model obtained from the repeated cross-validation in the training set to predict outcomes (that is, educational achievement, BMI and general cognitive ability) in the independent test set. To index prediction accuracy, we used the coefficient of determination, in the following referred to as $R^2_{test}$. Differences between mean-cv-$R^2_{train}$ and $R^2_{test}$ provide an index of out-of-sample error.

We used permutation to test the statistical significance of the difference in predictions between the MPS and the best single-score model. To test the null hypothesis of exchangeability of models, $H_0$: $_{MPS}R^2_{test} = _{best-single-score}R^2_{test}$, we compared the observed $_{diff}R^2_{test}$ ($_{MPS}R^2_{test} - _{best-single-score}R^2_{test}$) against an empirical null distribution of no difference in predictions between the MPS and the best single-score model. We tested the exchangeability of models by randomly selecting either the MPS or the best single-score model to generate predictions. We then calculated the difference in $R^2$ for two models with shuffled predictions. The process was repeated 100 000 times, generating an empirical null distribution of $_{diff}R^2$ under exchangeability of model predictions.

If the null hypothesis of no difference between models is true, it would not matter if we randomly exchange the model used for generating predictions. However, if the observed $_{diff}R^2_{test}$ value falls outside of those obtained when randomly exchanging models, this represents evidence against the null hypothesis of no difference in prediction between models. The statistical significance, as expressed in an empirical *P*-value, is calculated as the fraction of permutation values that are at least as extreme as the original $_{diff}R^2_{test}$ statistic observed in nonpermuted data.

## RESULTS

### MPS predictions

The MPS models showed better prediction in the independent test set than the best single-score models. The best single-score models were the large 2016 GWAS of years of education predicting 9.8% of the variance in educational achievement and 3.6% in general cognitive ability in the test set. For BMI, Obesity class 1 achieved the best single-score prediction, explaining 3.8% of the variance. (See Supplementary Table S2 for full single-score models results; see Supplementary Figure S2 for a visual overview of the single-score model results.) The MPS models explained 10.9% variance in educational achievement, 4.8% in cognitive ability and 5.4% in BMI in the test set. The improvement in variance explained compared with the best single-score models was 1.1% ($P = 4e-03$), 1.1% ($P = 2e-03$) and 1.6% ($P = 1e-04$), respectively.

Figures 1a–c show the polygenic predictors selected during training of the MPS models and their standardized coefficients. The ranking of predictors provides an index for their contributions to prediction. Analogous to conventional multiple regression, a standardized coefficient represents the contribution of the predictor to the outcome when adjusting for all other variables in the model.

The model predicting educational achievement retained 12 polygenic predictors (Figure 1a). Cognitive and socioeconomic polygenic scores took the top ranks. However, the psychiatric cross-disorder polygenic score, which aggregates genetic risk for bipolar disorder, schizophrenia, major depressive disorder, autism and attention deficit hyperactivity disorder, and the score for depressive symptoms in the general population were also retained by the model. The scores for Homeostasis Model Assessment of β-cell function, an index of β-cell function, and for coronary artery disease also contributed to prediction of educational achievement.

The MPS model predicting cognitive ability selected 10 polygenic scores during cross-validation (Figure 1b). The strongest contributions to prediction came from cognitive and socio-economic variables. Contributions from the psychiatric realm came from major depressive disorder, autism spectrum disorder and bipolar disorder, with the latter two having positive association with cognitive ability.

The MPS model predicting BMI retained 28 polygenic scores (Figure 1c). The top three strongest predictions came from obesity-related variables. Ranks four and five were taken by coronary artery disease and age at menarche (negative association). The sixth strongest predictor for children's BMI was the polygenic score based on the GWAS of mean caudate nucleus volume that plays a role in various non-motor functions including procedural and associative learning and inhibitory action control.[33–36] Other predictors included ulcerative colitis, leptin and neuroticism.

### Stratification by MPS

We examined the phenotypic values by quantile of the MPS distribution. Figures 2a–c plot the observed outcomes against the predictions by the MPS model in the test set. In general, the quantile results were roughly linear.

Figure 2a shows quantile results for mean exam grades. Individuals in the top 10% of the MPS distribution on average achieved an 'A' mean grade (across the three subjects Mathematics, English and Science), whereas individuals in the bottom 10% MPS distribution achieved a 'C' mean grade on average (top 10% mean = 9.74; bottom 10% mean = 8.33 (11 = A*, 10 = A, 9 = B, 8 = C, 7 = D, 6 = E, 5 = F, 4 = G, 0 = failed). Cohen's *d* was 1.20 (95% confidence interval 0.99–1.41) suggesting that 88% of the top 10% MPS group had a mean grade above that of the bottom 10% group, and there is an 80% probability that a person picked at

**Figure 1.** (**a**) Multi-polygenic score (MPS) model predicting educational achievement. Standardized coefficients of polygenic predictors selected by elastic net via repeated cross-validation in training set. Analogous to conventional multiple regression, a standardized coefficient represents the contribution of the predictor to the outcome when adjusting for all other variables in the model. The mean variance explained of the resampling distribution from the cross-validation was mean-cv-$R^2_{train} = 0.12$. The out-of-sample prediction of the model was $R^2_{test} = 0.109$. (**b**) MPS model predicting general cognitive ability. Standardized coefficients of polygenic predictors selected by elastic net via repeated cross-validation in training set. Analogous to conventional multiple regression, a standardized coefficient represents the contribution of the predictor to the outcome when adjusting for all other variables in the model. The mean variance explained of the resampling distribution from the cross-validation was mean-cv-$R^2_{train} = 0.051$. The out-of-sample prediction of the model was $R^2_{test} = 0.048$. (**c**) MPS model predicting body mass index (BMI). Standardized coefficients of polygenic predictors selected by elastic net via repeated cross-validation in training set. Analogous to conventional multiple regression, a standardized coefficient represents the contribution of the predictor to the outcome when adjusting for all other variables in the model. The mean variance explained of the resampling distribution from the cross-validation was mean-cv-$R^2_{train} = 0.074$. The out-of-sample prediction of the model was $R^2_{test} = 0.054$.

random from the top 10% MPS group will have a higher score than a person picked at random from the bottom 10% group.[37,38]

For cognitive ability, Figure 2b illustrates that individuals in the top 10% of the MPS distribution on average had a standardized cognitive ability score over 0.64 (95% confidence interval 0.40–0.89) s.d. higher than those in the bottom 10% MPS distribution. This means that 74% in the top 10% MPS group had mean ability score above that of the bottom 10% group, and that there is a

34

**Figure 2.** (**a**) Educational achievement by multi-polygenic score (MPS) deciles. Observed mean grade (across the three subjects Mathematics, English and Science) by deciles of the MPS predictions in the test set. Bars represent 95% confidence estimates. (**b**) General cognitive ability by MPS deciles. Observed mean standardized general cognitive ability by deciles of the MPS predictions in the test set. Bars represent 95% confidence estimates. (**c**) Body mass index (BMI) by MPS deciles. Observed mean standardized BMI (age and sex adjusted by external reference) by deciles of the MPS predictions in the test set. Bars represent 95% confidence estimates.

67% probability that a person picked at random from the top 10% MPS group will have a higher score than a person picked at random from the bottom 10% group.

For BMI, Figure 2c shows that children in the top 10% of the MPS distribution on average had a 0.80 (95% confidence interval 0.57–1.03) s.d. higher than those in the bottom 10% MPS distribution. Expressed differently, 79% of children in the top 10% MPS group had a mean ability score above that of the bottom 10% group,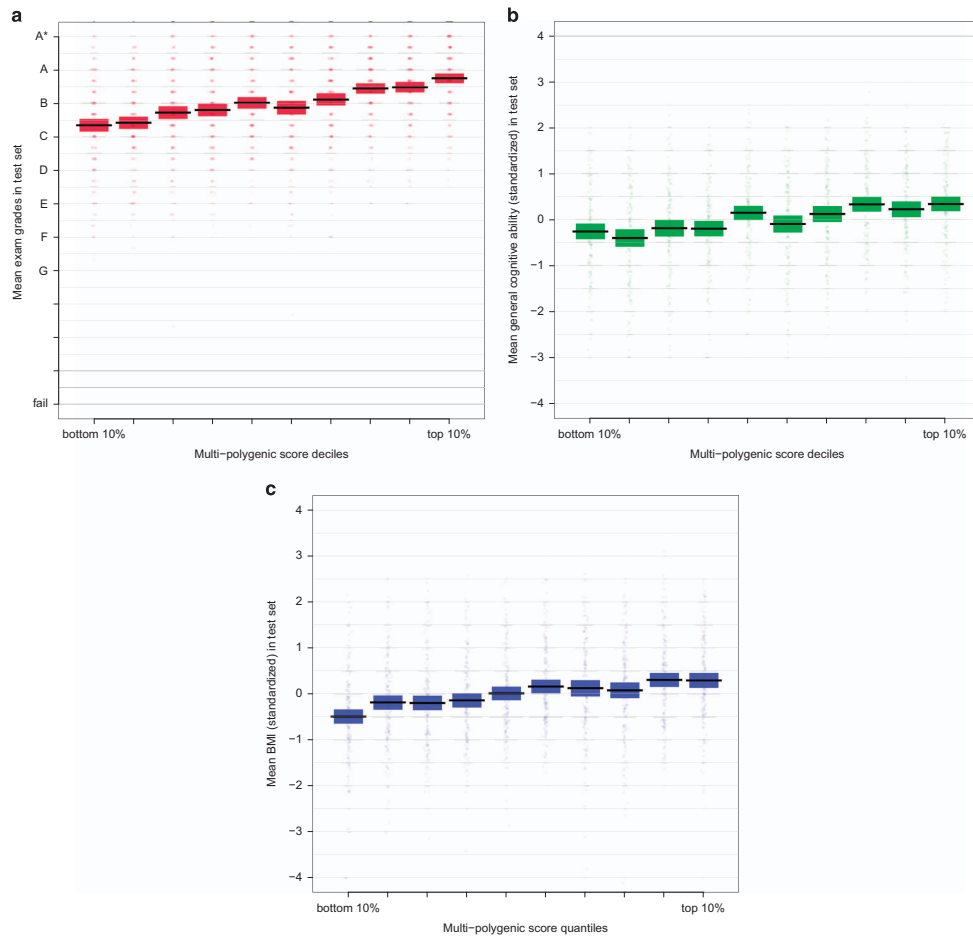 and that there is a 71% probability that a person picked at random from the top 10% MPS group will have a higher score than a person picked at random from the bottom 10% group.

**DISCUSSION**

We demonstrate that the MPS approach that combines summary-level GWAS data from multiple traits yields better individual-level phenotype prediction than single-score predictor models in independent test data.

The observation that a multitude of polygenic scores contribute to trait prediction in the MPS models highlights the complexity of the system being studied and the somewhat arbitrary way we divide it into phenotypic characteristics. We show that polygenic variation associated with traits other than the to-be-predicted outcome contributes to prediction. For instance, although there is a known association between ulcerative colitis and BMI,[39] genetic

variants associated with ulcerative colitis are not typically included in models estimating individuals' genetic risk for increased BMI.

The predictors selected and coefficients estimated by the MPS models in the current study can be used to generate individual-specific composite estimates of genetic propensities in other and smaller samples. For a more parsimonious replication, future research in other samples could construct a simple multiple regression model using the top five predictors selected by the current analyses. The predictive power of such an MPS model can then be compared with that of the best single-score model. More generally, in addition to the likely improvement in MPS prediction as more and larger GWASs are being published, the MPS approach has the potential to be applied to a wide range of outcomes and samples, including psychiatric and medical outcomes in case–control samples.

The predictive power of a polygenic score is not only a function of the genetic correlation between discovery and outcome trait, but also of the statistical power present in the discovery GWAS on which it is based (that is, variance explained and sample size).[5] The MPS approach exploits the fact that even GWASs of genetically distantly related traits might contribute predictive power if their power is superior to GWASs of more proximal traits. For instance, most likely because of its much greater sample size, the years of education polygenic score predicted general cognitive ability better than any of the polygenic scores based on GWASs directly measuring general cognitive ability.

Because predictive power of polygenic scores does not simply reflect the genetic correlation between discovery and target trait, but depends on the genetic architecture of both traits and sample size (especially of the discovery sample),[5,6,40] the MPS approach is not suited for investigating etiology. Other methods have been developed to that end. For instance, multivariate twin studies are appropriate for investigating trait etiology, or multi-trait GWAS meta-analysis aims to disentangle effects of correlated traits at the level of genetic variants.[15,16,41–45] In contrast, the declared aim of the MPS approach is to maximize trait prediction, without assumptions about the relationships among predictors.

The MPS approach will be useful whenever trait prediction is a priority. The primary reason for maximizing predictive power using the MPS approach is to predict phenotypes of individuals with as much accuracy as possible. Individual-specific genetic predictions will be useful in research with modest sample sizes to investigate developmental, multivariate and gene–environment interplay issues. Eventually, MPS models could be useful in both society and science to estimate genetic potential as well as risk in relation to all domains of functioning, including cognitive abilities and disabilities, personality and health and illness.

This predictive power will raise concerns about potential early, even prenatal, prediction. It is important to begin discussions that are informed by the empirical data because genotype-based trait prediction is moving towards the point of practical relevance. Although concerns are warranted, these might be outweighed by the benefits that could result from being able to predict problems and potential early and develop stratified preventions and interventions accordingly.

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## AUTHOR CONTRIBUTIONS

Conceived and designed the experiments: EK and RP. Analyzed data and processed and quality controlled genotype data: EK. Performed/supervised manual quality control and calling of genotype data: HP, SN and CJC. Wrote the paper: EK and RP. All authors contributed to and critically reviewed the manuscript.

## REFERENCES

1 Eicher JD, Landowski C, Stackhouse B, Sloan A, Chen W, Jensen N *et al.* GRASP v2.0: an update on the Genome-Wide Repository of Associations between SNPs and phenotypes. *Nucleic Acids Res* 2015; **43**: D799–D804.

2 Visscher PM, Brown MA, McCarthy MI, Yang J. Five years of GWAS discovery. *Am J Hum Genet* 2012; **90**: 7–24.

3 Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S *et al.* Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat Genet* 2014; **46**: 1173–1186.

4 Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* 2015; **518**: 197–206.

5 Dudbridge F. Power and predictive accuracy of polygenic risk scores. *PLoS Genet* 2013; **9**: e1003348.

6 Palla L, Dudbridge F. A fast method that uses polygenic scores to estimate the variance explained by genome-wide marker panels and the proportion of variants affecting a trait. *Am J Hum Genet* 2015; **97**: 250–259.

7 Euesden J, Lewis CM, O'Reilly PF. PRSice: polygenic risk score software. *Bioinformatics* 2014; **31**: 1466–1468, btu848.

8 Visscher PM, Yang J. A plethora of pleiotropy across complex traits. *Nat Genet* 2016; **48**: 707–708.

9 Bulik-Sullivan B, Finucane HK, Anttila V, Gusev A, Day FR, Loh P-R *et al.* An atlas of genetic correlations across human diseases and traits. *Nat Genet* 2015; **47**: 1236–1241.

10 Pickrell JK, Berisa T, Liu JZ, Ségurel L, Tung JY, Hinds DA. Detection and interpretation of shared genetic influences on 42 human traits. *Nat Genet* 2016; **48**: 709–717.

11 Krapohl E, Euesden J, Zabaneh D, Pingault J-B, Rimfeld K, von Stumm S *et al.* Phenome-wide analysis of genome-wide polygenic scores. *Mol Psychiatry* 2016; **21**: 1188–1193.

12 Solovieff N, Cotsapas C, Lee PH, Purcell SM, Smoller JW. Pleiotropy in complex traits: challenges and strategies. *Nat Rev Genet* 2013; **14**: 483–495.

13 Zou H, Hastie T. Regularization and variable selection via the elastic net. *J R Stat Soc Ser B Stat Methodol* 2005; **67**: 301–320.

14 Maier R, Moser G, Chen G-B, Ripke S, Cross-Disorder Working Group of the Psychiatric Genomics Consortium, Coryell W *et al.* Joint analysis of psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major depressive disorder. *Am J Hum Genet* 2015; **96**: 283–294.

15 Baselmans BM, Jansen R, Dongen J, van, Bao Y, Smart M, Kumari M *et al.* Multivariate genome-wide and integrated transcriptome and epigenome-wide analyses of the well-being spectrum. *bioRxiv* 2017; doi: 10.1101/115915.

16 Ferreira MAR, Purcell SM. A multivariate test of association. *Bioinformatics* 2009; **25**: 132–133.

17 Turley P, Walters RK, Maghzian O, Okbay A, Lee JJ, Fontana MA *et al.* MTAG: multi-trait analysis of GWAS. *bioRxiv* 2017; doi: 10.1101/118810.

18 Selzam S, Krapohl E, von Stumm S, O'Reilly PF, Rimfeld K, Kovas Y *et al.* Predicting educational achievement from DNA. *Mol Psychiatry* 2017; **22**: 267–272.

19 Krapohl E, Plomin R. Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs. *Mol Psychiatry* 2016; **21**: 437–443.

20 Kovas Y, Haworth CMA, Dale PS, Plomin R. The genetic and environmental origins of learning abilities and disabilities in the early school years. *Monogr Soc Res Child Dev* 2007; **72**, vii 1–144.

21 McCarthy S, Das S, Kretzschmar W, Durbin R, Abecasis G, Marchini J. A reference panel of 64,976 haplotypes for genotype imputation. *bioRxiv* 2015; **48**: 1279–1283, 35170.

22 Zheng J, Erzurumluoglu AM, Elsworth BL, Kemp JP, Howe L, Haycock PC *et al.* 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* 2017; **33**: 272–279.

23 Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al.* LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat Genet* 2015; **47**: 291–295.

24 Vilhjálmsson BJ, Yang J, Finucane HK, Gusev A, Lindström S, Ripke S *et al.* Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *Am J Hum Genet* 2015; **97**: 576–592.

25 Kuhn M. Building predictive models in R using the caret package. *J Stat Softw* 2008; **28**: 1–26.

26 Tibshirani R. Regression shrinkage and selection via the Lasso. *J R Stat Soc Ser B Methodol* 1996; **58**: 267–288.

27 Zhou D-X. On grouping effect of elastic net. *Stat Probab Lett* 2013; **83**: 2108–2112.

28 Friedman J, Hastie T, Tibshirani R. Regularization paths for generalized linear models via coordinate descent. *J Stat Softw* 2010; **33**: 1–22.

29 Kohavi R. *A study of cross-validation and bootstrap for accuracy estimation and model selection.* In: Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2. Morgan Kaufmann Publishers Inc.: San Francisco, CA, USA, 1995 (cited 14 December 2016), pp 1137–1143. Available from http://dl.acm.org/citation.cfm?id = 1643031.1643047.

30 Efron B. Estimating the error rate of a prediction rule: improvement on cross-validation. *J Am Stat Assoc* 1983; **78**: 316–331.

31 Efron B, Tibshirani R. Improvements on cross-validation: the 632+ bootstrap method. *J Am Stat Assoc* 1997; **92**: 548–560.

32 Kim J-H. Estimating classification error rate: repeated cross-validation, repeated hold-out and bootstrap. *Comput Stat Data Anal* 2009; **53**: 3735–3745.

33 Malenka RC, Nestler E, Hyman S, Sydor A, Brown R. *Molecular Neuropharmacology: A Foundation for Clinical Neuroscience.* McGraw Hill Medical Book: New York, 2009.

34 Aron AR, Schlaghecken F, Fletcher PC, Bullmore ET, Eimer M, Barker R *et al.* Inhibition of subliminally primed responses is mediated by the caudate and thalamus: evidence from functional MRI and Huntington's disease. *Brain* 2003; **126**: 713–723.

35 Jahanshahi M, Obeso I, Rothwell JC, Obeso JA. A fronto-striato-subthalamic-pallidal network for goal-directed and habitual inhibition. *Nat Rev Neurosci* 2015; **16**: 719–732.

36 Seger CA, Cincotta CM. The roles of the caudate nucleus in human classification learning. *J Neurosci* 2005; **25**: 2941–2951.

37 Ruscio J. A probability-based measure of effect size: robustness to base rates and other factors. *Psychol Methods* 2008; **13**: 19–30.

38 Cohen J. *Statistical Power Analysis for the Behavioral Sciences,* 2nd edn. Academic Press: New York, NY, US, 1988.

39 Dong J, Chen Y, Tang Y, Xu F, Yu C, Li Y *et al.* Body mass index is associated with inflammatory bowel disease: a systematic review and meta-analysis. *PLoS ONE* 2015; **10**: e0144872.

40 Wray NR, Lee SH, Mehta D, Vinkhuyzen AAE, Dudbridge F, Middeldorp CM. Research review: polygenic methods and their application to psychiatric traits. *J Child Psychol Psychiatry* 2014; **55**: 1068–1087.

41 Majumdar A, Haldar T, Bhattacharya S, Witte J. An efficient Bayesian meta-analysis approach for studying cross-phenotype genetic associations. bioRxiv. 2017; doi: 10.1101/101543.

42 Zhu X, Feng T, Tayo BO, Liang J, Young JH, Franceschini N *et al.* Meta-analysis of correlated traits via summary statistics from GWASs with an application in hypertension. *Am J Hum Genet* 2015; **96**: 21–36.

43 Bhattacharjee S, Rajaraman P, Jacobs KB, Wheeler WA, Melin BS, Hartge P *et al.* A subset-based approach improves power and interpretation for the combined analysis of genetic association studies of heterogeneous traits. *Am J Hum Genet* 2012; **90**: 821–835.

44 Shim H, Chasman DI, Smith JD, Mora S, Ridker PM, Nickerson DA *et al.* A multivariate genome-wide association analysis of 10 LDL subfractions, and their response to statin treatment, in 1868 Caucasians. *PLOS ONE* 2015; **10**: e0120758.

45 Bolormaa S, Pryce JE, Reverter A, Zhang Y, Barendse W, Kemper K *et al.* A multi-trait, meta-analysis for detecting pleiotropic polymorphisms for stature, fatness and reproduction in beef cattle. *PLOS Genet* 2014; **10**: e1004198.

Supplementary Information accompanies the paper on the Molecular Psychiatry website (http://www.nature.com/mp)

# 5 Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs

Krapohl, E. and Plomin, R. (2016). Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs. *Molecular Psychiatry*, 21(3):437–443

Supplementary material: `http://www.nature.com/mp/journal/v21/n3/suppinfo/mp20152s1.html`

## ORIGINAL ARTICLE

# Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs

E Krapohl and R Plomin

One of the best predictors of children's educational achievement is their family's socioeconomic status (SES), but the degree to which this association is genetically mediated remains unclear. For 3000 UK-representative unrelated children we found that genome-wide single-nucleotide polymorphisms could explain a third of the variance of scores on an age-16 UK national examination of educational achievement and half of the correlation between their scores and family SES. Moreover, genome-wide polygenic scores based on a previously published genome-wide association meta-analysis of total number of years in education accounted for ~ 3.0% variance in educational achievement and ~ 2.5% in family SES. This study provides the first molecular evidence for substantial genetic influence on differences in children's educational achievement and its association with family SES.

*Molecular Psychiatry* (2016) **21,** 437–443; doi:10.1038/mp.2015.2; published online 10 March 2015

## INTRODUCTION

After health care, education is society's largest and most expensive environmental intervention, consuming > 6% of gross domestic product in OECD (Organization for Economic Co-operation and Development) countries.[1] Understanding the etiology and correlates of differences between children in what they take away from their education is important because their educational achievement directly determines admission to further education and employability and also predicts a wide range of health outcomes.[1–3] Pedigree-based methods, primarily twin studies comparing the similarity of identical and nonidentical twins, have consistently suggested substantial genetic influence on differences between children in their educational achievement.[4–10] It is now possible to use DNA-based methods to estimate genetic influence on variance in large samples of unrelated individuals.[11,12] No DNA-based estimates of genetic influence have as yet been reported for children's educational achievement, although evidence has been reported for the rough proxy of total number of years in education in adults.[13–16] This study used children's genotypes to estimate genetic influences on variance in educational achievement and its covariance with family socioeconomic status (SES).

Here we report the first investigation of genetic influence on the variance of children's educational achievement using DNA alone. The same DNA-based methods can also be used to estimate genetic influence on the covariance between traits.[17] This enabled us to investigate possible genetic mediation of the best predictor of children's educational achievement, their family's SES.[18,19] This correlation is often interpreted causally as family SES causing differences in children's educational achievement.[20] However, it remains unclear whether and to what extent the association between family SES and children's educational achievement is genetically mediated, because twin and family

research is limited to studying phenotypes that can vary within a family. Key aspects of children's environment such as poverty, parental education and neighborhood cannot be investigated using the twin method because it is methodologically impossible to decompose variance in phenotypes shared within twin pairs.

The DNA-based technique, genome-wide complex trait analysis (GCTA),[11] fits the effects of genome-wide single-nucleotide polymorphisms (SNPs) as random effects in a mixed linear model to estimate variance or covariance captured by all SNPs simultaneously. Contrary to traditional family-based methods that estimate the genetic contribution to phenotypic variation or covariation by known kinship coefficients, GCTA relies on empirical genetic resemblance established from identity by state inferred from genome-wide SNP similarity of 'unrelated' individuals.

Because GCTA is based on unrelated individuals, it enables the decomposition of variance of phenotypes such as family SES that are the same for members of a family and therefore cannot be decomposed in analyses such as the twin method that rely on within-family differences. Another difference between the two methods is that, unlike the twin method, GCTA is limited to estimating additive genetic effects for the SNPs on the genome-wide DNA array or other DNA variants in linkage disequilibrium with the measured SNPs, which until recently have been common SNPs. Thus, GCTA will underestimate genetic influence to the extent that nonadditive effects or rare variants contribute importantly to heritability. This limitation of GCTA to additive effects of common SNPs is the same limitation of genome-wide association (GWA) studies that attempt to identify specific SNPs associated with a trait. GCTA is directly comparable to GWA results because both rely on the same experimental design using the same genetic signal;[21] GCTA provides an upper-limit estimate of the genetic effects that can be identified by GWA.

King's College London, MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology & Neuroscience, London, UK. Correspondence: E Krapohl, King's College London, Social, Genetic and Developmental Psychiatry Research Centre, Institute of Psychiatry, Psychology & Neuroscience, De Crespigny Park, Denmark Hill, London SE5 8AF, UK.
E-mail: eva.krapohl@kcl.ac.uk

GWA attempts aimed at identifying individually significant SNPs have generally captured only extremely small fractions of genetic variance of complex traits, the so-called *missing heritability* problem.[22] However, evidence has been accumulating that significant portions of phenotypic variation can be explained by the ensemble of markers not achieving genome-wide significance.[23] Markers are identified from GWAs using an initial discovery sample to construct a genome-wide polygenic score (GPS) in an independent replication sample by calculating the effect-size-weighted sum of trait-associated alleles for each individual. An aggregate GPS score can be used to assess genetic influence on trait variation.

As they are tapping into the same genetic signal, GPS based on GWA results and GCTA can be applied to the same data sets, with both estimating the polygenic contribution to trait variance or a shared polygenic covariance between traits captured by the additive effects of common SNPs. We therefore employ a two-method approach using GCTA and GPS to explore the genetic influence on the variance of children's educational achievement and on the covariance between family SES and children's educational achievement. Our study had four objectives:

(1) To estimate, for the first time using DNA data, genetic influences on children's educational achievement on an age-16 UK national examination of educational achievement using genome-wide genotypes from >3000 conventionally unrelated children. Specifically, we conduct GCTA[11] to quantify pairwise genomic similarity between each pair of individuals across millions of SNPs throughout the genome in order to estimate the proportion of phenotypic variation in children's educational achievement captured by all SNPs simultaneously.

(2) To investigate genetic mediation of the phenotypic correlation between family SES and children's educational achievement, we conduct bivariate GCTA to estimate the proportion of phenotypic covariation between children's family SES and children's educational achievement captured by children's genotypes.

(3) To create a GPS based on the results of a large GWA study on adults' total years of schooling[13] and investigate its association with variance in children's educational achievement and their family SES.

(4) To examine the role of general cognitive ability (intelligence) in the genetic nexus between children's educational achievement and their family SES. Molecular evidence as well as twin studies have shown that cognitive ability is heritable and accounts for substantial portion of genetic variance in educational achievement.[7,24–26] In addition, recent molecular evidence from the present sample of unrelated individuals showed high genetic correlation between family SES and children's intelligence at age 7 and 12 years.[27] Based on this evidence, it is important to address the question to what extent the genetic link between family SES and children's educational achievement is mediated by intelligence. For this reason, we perform GCTA mediation analyses to test for a direct genetic link between family SES and children's educational achievement independent of cognitive ability. Complementarily, we test whether the GPS of adults' total years of schooling explains variance in children's educational achievement independently of cognitive ability.

Our findings provide the first molecular evidence for substantial genetic influence on variation in children's educational achievement and its association with family SES. We further show that children's intelligence accounts for one third of this SNP link between family SES and children's educational achievement. In addition, we demonstrate that a GPS based on years of education in adulthood discovered in an independent large GWA meta-analysis[13] significantly explains variance in children's educational achievement in our sample, even after controlling for intelligence.

## MATERIALS AND METHODS

### Sample and genotyping

The sample was drawn from the Twins Early Development Study (TEDS), a multivariate longitudinal study that recruited over 11 000 twin pairs born in England and Wales in 1994, 1995 and 1996.[28,29] TEDS has been shown to be representative of the UK population.[30] Supplementary Table 2 shows that the genotyped subsample of TEDS is representative of UK census data from first contact through age 16 years.

The project received approval from the Institute of Psychiatry ethics committee (05/Q0706/228) and parental consent was obtained before data collection.

DNA data were available for 3747 children whose first language was English and had no major medical or psychiatric problems. From that sample, 3665 DNA samples were successfully hybridized to Affymetrix GeneChip 6.0 SNP genotyping arrays (Affymetrix, Santa Clara, CA, USA) using standard experimental protocols as part of the WTCCC2 project (for details see Trzaskowski et al.).[31] In addition to nearly 700 000 genotyped SNPs, more than one million other SNPs were imputed from HapMap 2, 3 and WTCCC controls using IMPUTE v.2 software.[32] A total of 3152 DNA samples (1446 males and 1706 females) survived quality control criteria for ancestry, heterozygosity, relatedness and hybridization intensity outliers. To control for ancestral stratification, we performed principal component analyses on a subset of 100 000 quality-controlled SNPs after removing SNPs in linkage disequilibrium ($r^2 > 0.2$).[33] Using the Tracy–Widom test,[34] we identified 8 axes with $P < 0.05$ that were used as covariates in GCTA and polygenic score analyses.

### Measures

*Educational achievement.* Educational achievement was operationalized as performance on the standardized UK-wide examination, the General Certificate of Secondary Education (GCSE), taken by almost all (>99%) pupils at the end of compulsory education at typically at the age of 16 years. English, mathematics and science are compulsory subjects. Five or more GCSEs with grades A*–C are required for further education, including GCSE English and GCSE mathematics. The joint performance on these three compulsory subjects determines admission to further education and employability.

The data for the present study were collected by questionnaires sent by mail and by telephone interview of parents and twins themselves. After completed forms were received from the families, the grades were coded from 11 (the highest grade: A*) to 4 (the lowest pass grade: G); no information about failed results was available. For 1729 individuals, self- and parent-reported GCSE results were verified using data obtained from the UK National Pupil Database,[35] yielding correlations of 0.99 for mathematics, 0.98 for English and 0.96 for science.

The GCSE measure for the present analyses was the mean grade of the three compulsory core subjects, mathematics, English (mean grade of 'English Language' and 'English Literature'), and science (mean of any science subjects taken), requiring at least two measures to be nonmissing. Scores on the three compulsory core subjects were highly correlated (0.65–0.81).

*Intelligence (IQ).* Individuals were assessed at the ages of 2, 3, 4, 7, 9, 10, 12, 14, and 16 years on general cognitive ability using a battery of parent-administered and phone- and web-based tests. At ages 2, 3, and 4, tests were parent-administered and validated against standard tests administered by a trained tester. At age 7, tests were administered over the phone; at age 9, parents administered the tests; and at the ages 10 – 16, tests were web based. At each testing age, individuals completed at least two ability tests that assessed verbal and nonverbal intelligence. Psychometric properties of the tests have been described in detail elsewhere,[36] with the exception of the measurements used at age 16 years, where subjects completed a web-based adaptation of Raven's Standard and Advanced Progressive Matrices and the Mill-Hill Vocabulary Scale.[37–39]

For each composite measure at each of the nine ages, scores were regressed on sex and age, outliers above or below 3 s.d. from the mean were excluded and the standardized residuals were quantile normalized. Subsequently, a mean composite scale was created as the mean across the nine ages, performing mean-imputation for missing measurement occasions to avoid list-wise deletion.

*Family SES.* Converging evidence suggests that a composite of variables including parental education and occupation represents SES better than any single indicator.[18] To index family SES, we combined parental

education and occupation assessed when children were aged 2, 7 and 16 years. At age 2 years, SES was constructed as the mean of mother's and father's highest education level, mother's and father's occupation assessed by the Standard Occupational Classification 2000,[40,41] and maternal age at birth of eldest child. The SES composite when children were age 7 years was created similarly but without the variable of age of mother at birth of eldest child. At age 16 years, SES was composed as the mean of household income, maternal and paternal education level and maternal and paternal occupation. Mean composites were standardized and quantile normalized. The correlations between these three SES estimates ranged from 0.70 to 0.77. To increase reliability and maximize sample size, the final measure of family SES for this study was created as the mean composite score of parental SES reported when children were aged 2, 7, and 16 years, performing mean-imputation for missing data points.

### Statistical analyses

*GCTA.* The GCTA model decomposes the trait variance into an additive genetic component ($_G$) captured by the available SNPs (and correlated markers in linkage disequilibrium with the genotyped SNPs) and a residual component containing all nonadditive genetic variance, interaction effects, environmental factors, error variance and additive genetic variance that is not tagged by the sampled SNPs. Hence, the GCTA model estimates lower-bound additive genetic variance for both phenotypes ($V_G^{GCSE}$, $V_G^{SES}$); and the correlation between the additive genetic components ($\rho_G$). The $\rho_G$ is not biased in the same way $V_G$ is. This is because the estimate of genetic correlation is a function of the ratio between SNP-tagged covariance and SNP-tagged variance that are biased to the same extent (that is, the estimates are subject to the same imperfect linkage disequilibrium between causal variants and genotyped SNPs) and hence cancel each other out.[42]

Using genome-wide SNP data, we estimate genetic variation and covariation from a representative sample of 3000 unrelated children. Our estimates were obtained by restricted maximum likelihood using the published algorithm for GCTA.[11] GCTA estimates the proportion of phenotypic variance of a trait tagged by sampled SNPs by fitting the polygenic effects of all SNPs simultaneously as random effects in a mixed linear model using a restricted maximum likelihood function. The so-called genetic relatedness matrix holds the mean pairwise genomic similarity (weighted by allele frequency) between all pairs of individuals in the sample across all SNPs. The variance tagged by all SNPs is estimated to be >0 when genetically more similar individuals are phenotypically more similar. The bivariate extension of the model relates the pairwise genetic similarity matrix to a phenotypic covariance matrix between traits (here family SES and educational achievement).[17] To prevent confounding of the SNP estimate by shared environment effects and the effects of causal variants that are not tagged by the SNPs, cryptic relatedness was removed from the analyses. This default procedure eliminates one individual from a pair whose genetic similarity is 0.025 or greater; a coefficient that approximates at least fifth-degree relatives. The removal of close relatives ensures that estimates reflect the tagging of causal variants through population linkage disequilibrium. This criterion removed seven individuals from the analyses. Analyses were executed using GCTA[11] and R software.[43]

The present sample size of ~3000 yields 80% power to detect a GCTA heritability estimate of 30% (α = 0.05) and genetic correlation estimate of 0.6 (α = 0.05; $V_G^1$ = 0.20; $V_G^2$: 0.30; $r_{Ph}$ = 0.50).

*Polygenic scores.* We created polygenic scores from genome-wide data of over 3000 unrelated children using GWA results for total years of schooling from an independent discovery sample.[13] The same quality control criteria as for the GCTA analyses were applied to the data. Polygenic risk scores were constructed using the *P*-values and β-weights from the recent large (*N* = 126 559) GWA based on years of education.[6] Quality-controlled SNPs were pruned for linkage disequilibrium based on *P*-value informed clumping in PLINK,[44] using $R^2$ = 0.25 cutoff within a 200-kb window. We removed the major histocompatibility complex region of the genome because of its complex linkage disequilibrium structure. 144 890 SNPs survived linkage disequilibrium pruning. For each individual, multiple polygenic scores were generated using the PLINK score option based on the top SNPs from the GWA analysis of educational attainment for varying significance thresholds (from 0.01 to 0.50). Numbers of SNPs per threshold are summarized in Supplementary Table 3. The scores were calculated as the sum across SNPs of

the number of reference alleles for each SNP multiplied by the effect size (β-coefficient) derived from the GWA analysis of years of education.

Polygenic scores were tested for association with the same quantitative measures used in the GCTA analyses (family SES, educational achievement (GCSE), intelligence and educational achievement controlled for intelligence) in linear regressions. These analyses were corrected for the first eight ancestry-informative principal components by entering them as covariates into the regression models. Analyses were performed in PLINK and R.

### RESULTS

Phenotypically, children's educational achievement correlated 0.50 (0.02 s.e.) with their family SES. Both variables also correlated with intelligence: 0.55 (0.02 s.e.) for educational achievement and 0.38 (0.02 s.e.) for family SES (Supplementary Table 1).

### Bivariate GCTA

Bivariate GCTA showed that the estimated proportion of variance tagged by the sampled SNPs was 0.31 (0.12 s.e.) in educational achievement, and 0.20 (0.11 s.e.) in family SES (Figure 1). The genetic correlation, indicating the extent to which the same SNPs are associated with family SES and children's educational achievement, was near unity ($r_G$ = 1.02 (0.25 s.e.)).

Based on the genetic correlation between the two traits and the genetic contribution to variance of each trait respectively, GCTA estimates the genetic contribution to the phenotypic correlation between the two traits: $C_{(G)} = r_{1,2}$ $_{(G)} \sqrt{(V_1$ $_{(G)} \times V_2$ $_{(G)})}$, applied to the data: $0.25 = 1.02 \times \sqrt{(0.31 \times 0.20)}$. Hence, GCTA estimated the genetic contribution to the phenotypic correlation between family SES and children's educational achievement as 0.25 (0.09 s.e.), indicating that the proportion of the observed correlation tagged by the additive effects of available SNPs was 50% (that is, 0.25/0.50; Figure 1). This suggests approximately half of the phenotypic correlation between children's family SES and their educational achievement was mediated genetically.

*Mediation analyses.* To test whether intelligence mediates the observed association between family SES and children's educational achievement, we statistically controlled for intelligence by regressing GCSE on intelligence and entering the resulting standardized residuals into the bivariate GCTA model with family SES. When controlling for variance explained by children's intelligence, which yielded a univariate GCTA estimate of 0.38 (0.11 s.e.) (data not shown), the phenotypic correlation between family SES and children's educational achievement was reduced from 0.50 to 0.37 (0.02 s.e.). The GCTA estimate of the genetic covariation between family SES and children's educational achievement dropped from 0.25 (0.09 s.e.) to 0.17 (0.09 s.e.). Mirroring the mediation observed at the phenotypic level, this suggests that one-third of the SNPs tagging variation in family SES and children's educational achievement also captured individual differences in intelligence, implying two-thirds of the SNPs linking family SES and children's educational achievement were independent of intelligence.

### Polygenic score analyses

Polygenic score analysis is designed to test whether SNPs that do not reach genome-wide significance in a discovery GWA are nonetheless significantly associated in aggregate with a trait in an independent sample. In the same sample of 3152 unrelated individuals, we created polygenic scores with varying numbers of SNPs (see Materials and methods) based on a large meta-analytic GWA study (*N* = 126 599) of years of education.[13] Figure 2 displays the results of multiple linear regression analyses showing that the polygenic scores accounted for ~3.0% variance in educational

**Figure 1.** Bivariate genome-wide complex trait analysis (GCTA) of family socioeconomic status (SES) and children's educational achievement (General Certificate of Secondary Education (GCSE)). (**a**) Proportion of phenotypic trait variance tagged by the sampled SNPs in GCSE and family SES, respectively. (**b**) Covariance between family SES and GCSE captured by SNPs, without controlling for intelligence (left bar) and when controlling for intelligence (GCSE.IQ) (right bar). The length of the bar indicates the total phenotypic correlation between SES and GCSE. Solid black lines indicate standard errors.

achievement (GCSE), ~2.5% in family SES and ~1.0% in intelligence. All $P$-values were $\leq 3.79^{-07}$. Notably, the effect size for GCSE remained substantial (~2.0%) and significant ($P \leq 2.27^{-06}$) when statistically controlling for intelligence.

## DISCUSSION

This study provides the first molecular evidence for substantial genetic influence on differences in children's educational achievement at the end of compulsory education in the United Kingdom and its association with family SES. Our GCTA results show that SNPs that are associated with both family SES and GCSE scores account for about half of the phenotypic correlation between SES and GCSE. Mediation analysis suggests that about one-third of this genetic effect also extends to children's intelligence, but two-thirds of the genetic association between family SES and GCSE scores is independent of intelligence. In GPS analysis, we show that SNPs associated with total years of education in adulthood discovered by an independent large GWA meta-analysis[13] explain up to 3% of the variance in children's educational achievement in our sample, and up to 2% of the variance after controlling for intelligence.

The GCTA heritability estimate of 31% for children's performance on a UK national examination at the end of compulsory education corroborates the vast literature of traditional family-based methods, mostly the twin method, showing that variation in children's educational achievement is under substantial genetic influence,[4,5,7–9,45,46] with heritability estimates converging at ~50%. This commonly observed discrepancy in phenotypic

variance explained by pedigree-based methods (that is, twin and family) and population-based methods (that is, GCTA) occurs because GCTA only captures genetic variance contributed by additive effects of common SNPs that are in sufficient linkage disequilibrium with the causal DNA variants.[47]

Our GCTA heritability estimate of 20% for family SES tagged by children's genotypes is very similar to GCTA heritability estimates of years of education in adulthood and socioeconomic measures tagged by adults' genotypes themselves in previous studies.[13–15] This is remarkable as children's genotypes are only a proxy for their parents' genotypes. In other words, GCTA effects on family SES estimated from children's DNA only reflect the extent to which children inherit parental characteristics associated with the family SES created by the parents. One such factor is intelligence, and we find that children's intelligence accounts for about one-third of the GCTA association between family SES and children's educational achievement. However, it is interesting that two-thirds of the GCTA association is *not* accounted for by children's intelligence. This finding of intelligence-independent shared genetic variance between family SES and children's educational achievement suggests that differences in educational achievement at the end of compulsory education and the level of education and occupation attained in adulthood are not merely the manifestation of differences in intelligence. This is in line with twin research that suggests that the heritability of educational achievement reflects many genetically influenced traits such as personality and self-efficacy, not just intelligence.[48]

42

**Polygenic scores for education predict GCSE, family SES, and intelligence**



**Figure 2.** Genome-wide polygenic scores (GPS) for years of schooling in adults (Rietveld et al.[13]) predict variance ($R^2$) in children's educational achievement (General Certificate of Secondary Education (GCSE)), family socioeconomic status (SES), intelligence and educational achievement after controlling for intelligence (GCSE.IQ). GPS were created using different significance thresholds for inclusion of variants for years of education, ranging from $P = 0.01$ to $0.50$, indicated by heat colors. The uncorrected $P$-values above each bar indicate the statistical significance of the observed association between the GPS and the respective trait.

The polygenic nature of behavioral traits poses a statistical challenge as enormous sample sizes are needed to identify genome-wide significant single DNA variants.[23] Therefore, genome-wide methods, such as GCTA and GPS analysis, that aggregate genetic effects across a multitude of markers have the assumption of polygenicity at their core and provide powerful approaches for exploring genetic influences on traits and shared between traits.

A GPS based on markers associated with years of education in adulthood in an independent discovery sample was significantly associated with children's educational achievement in our sample. Replicating results from polygenic score analyses of a recent Dutch study,[49] this shows that the shared polygenic link between children's educational achievement and adult measures of education even holds when limited to education-associated SNPs identified in an independent sample of adults. We further demonstrate that this polygenic link persists independently of children's cognitive ability, and that the educational attainment GPS of children's genotypes explains variance in their parents' socioeconomic status. The predictive power of GPS analysis in our independent sample illustrates that adequately powered GWA studies can identify replicable genetic associations with behavioral traits. Although the current GPS accounts for only a small amount of phenotypic variance, as prediction improves, GPS can identify profiles of genetic risk and protective factors for unrelated individuals, which will enable more powerful prediction models that combine genetic and nongenetic factors. Polygenic

predictors might also facilitate research on the causal pathways underlying these genetic predictors.[21,22,50]

The results need to be interpreted in the context of three main important methodological limitations. First, a specific limitation of this study is its modest statistical power in the GCTA analyses (see Materials and methods). The GPS analyses were sufficiently powered to identify trait-associated variance at high statistical significance, but were limited by the power of the discovery GWAS to detect the small effect sizes of single variants across the genome.[21,23] A second, general limitation is the allelic spectrum covered by the current DNA microarrays, such as the Affymetrix 6.0 GeneChip used in our study, that is restricted to common variants. Research has begun exploring the relative contribution of common and rare variants to variation of psychiatric traits (see, for example, Gaugler et al.[51] and Yang et al.[52]). Future studies with greater statistical power may explore the relative contribution of common and rare variants to trait variation of educational achievement and associated phenotypes. Third, both GCTA and GWAS, on which GPS analysis relies, are limited to detecting additive genetic variation that is captured by the sampled SNPs, which are typically common SNPs with minor allele frequencies > 0.05. Hence, GCTA heritability provides a lower-bound narrow-sense heritability estimate and represents the upper limit for detection of SNP associations in GWA studies and thus for GPS analysis. Generally, these limitations imply a substantial underestimation of 'true heritability' in the present analyses.

The present analyses demonstrate the ability of DNA-based methods to explore the genetic architecture of extended

43

phenotypes such as family SES that cannot be detected by traditional variance/covariance estimation methods that rely on known kinship relatedness. Quantitative DNA-based methods, which rely on empirically established pairwise genomic similarity among traditionally unrelated individuals, can supplement and extend family-based methods and thereby facilitate the move from behavioral genetics to behavioral genomics.

Importantly, no directionality or causality can be inferred from the present results. Heritability indexes the proportion of trait variance attributable to genetic effects in a particular population at a particular time.[53] Finding evidence for heritability of a trait or co-heritability of two traits does not imply resistance to environmental factors as genetic effects are dynamic and subject to developmental and environmental change.[54] Research on how the heritability of educational achievement differs across development and across context suggests that genetic influences on these phenotypes are maximized by environmental opportunity.[54–56] Differences in individuals' exposure to environments are not random. Genotype–environment correlation refers to the empirical observation that individuals experience different environments as a systematic function of their genotypes.[56–61] Genetic effects on phenotypes may be mediated through developmental or socio-contextual processes.

Our results also contribute to the extensive debate about meritocracy and social mobility[62] that has largely ignored the fact that parents and their offspring are genetically related. Usually a lower correlation between parental and offspring SES is seen as an index of social mobility.[63] However, considering genetics, we know that removing environmental sources of variation will not remove genetically driven resemblance between parents and offspring. To the contrary, as environmental differences diminish, individual differences that remain will to a larger proportion be due to genetic differences; that is, heritability would increase, which has also been demonstrated empirically.[55] That way, heritability could be seen as an index of social mobility.

No necessary policy implications arise from finding heritability of educational achievement and its link with family SES. However, consideration of empirical evidence will lead to better-informed policy decisions. Specifically, analogous to the long-established model of evidence-based medicine, we believe that evidence-based education facilitated by a dialog between scientists and policy makers will be beneficial to education of all children and can also benefit schools, teachers, and society at large.[64]

In summary, our GCTA results show a substantial contribution of common SNPs to variation in children's educational achievement and its association with family SES. This is further substantiated by the GPS analyses, revealing significant sharing of genetic variants between children's educational achievement and total years of education in adulthood. Together, these findings provide converging evidence for substantial genetic influence on differences in children's educational achievement and genetic links with family SES. Our findings add weight to the view that genetic variation plays an important, but not exclusive, role in educational inequalities and social mobility, which is at variance with views, that still prevail in some quarters, that these are solely the product of social forces and environmental inequalities.

## CONFLICT OF INTEREST

The authors declare no conflict of interest.

## AUTHOR CONTRIBUTIONS

RP directs and received funding for the Twins Early Development Study (TEDS). EK conceived of the present study and analyzed the data. EK and RP wrote the manuscript.

## REFERENCES

1 OECD. *Education at a Glance 2013*. Organisation for Economic Co-operation and Development, 2013 Available at http://www.oecd-ilibrary.org/content/book/eag_highlights-2013-en.
2 Morris JN, Blane DB, White IR. Levels of mortality, education, and social conditions in the 107 local education authority areas of England. *J Epidemiol Community Health* 1996; **50**: 15–17.
3 White IR, Blane D, Morris JN, Mourouga P. Educational attainment, deprivation-affluence and self reported health in Britain: a cross sectional study. *J Epidemiol Community Health* 1999; **53**: 535–541.
4 Haworth CMA, Plomin R. Quantitative genetics in the era of molecular genetics: learning abilities and disabilities as an example. *J Am Acad Child Adolesc Psychiatry* 2010; **49**: 783–793.
5 Shakeshaft NG, Trzaskowski M, McMillan A, Rimfeld K, Krapohl E, Haworth CM et al. Strong genetic influence on a UK nationwide test of educational achievement at the end of compulsory education at age 16. *PLoS One* 2013; **8**: e80341.
6 Johnson W, Deary IJ, Iacono WG. Genetic and environmental transactions underlying educational attainment. *Intelligence* 2009; **37**: 466–478.
7 Calvin CM, Deary IJ, Webbink D, Smith P, Fernandes C, Lee SH et al. Multivariate genetic analyses of cognition and academic achievement from two population samples of 174,000 and 166,000 school children. *Behav Genet* 2012; **42**: 699–710.
8 Martin NG, Martin PG. The inheritance of scholastic abilities in a sample of twins I. Ascertainment of the sample and diagnosis of zygosity. *Ann Hum Genet* 1975; **39**: 213–218.
9 Gill CE, Jardine R, Martin NG. Further evidence for genetic influences on educational achievement. *Br J Educ Psychol* 1985; **55**: 240–250.
10 Haworth CMA, Plomin R. Genetics and education: Toward a genetically sensitive classroom. In: Harris KR, Graham S, Urdan T, McCormick CB, Sinatra GM, Sweller J (eds). *APA educational psychology handbook, Vol 1: Theories, constructs, and critical issues*. American Psychological Association: Washington, DC, USA, 2012, pp 529–559.
11 Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet* 2011; **88**: 76–82.
12 Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR et al. Common SNPs explain a large proportion of the heritability for human height. *Nat Genet* 2010; **42**: 565–569.
13 Rietveld CA, Medland SE, Derringer J, Yang J, Esko T, Martin NW et al. GWAS of 126,559 individuals identifies genetic variants associated with educational attainment. *Science* 2013; **340**: 1467–1471.
14 Marioni RE, Davies G, Hayward C, Liewald D, Kerr SM, Campbell A et al. Molecular genetic contributions to socioeconomic status and intelligence. *Intelligence* 2014; **44**: 26–32.
15 Benjamin DJ, Cesarini D, van der Loos MJ, Dawes CT, Koellinger PD, Magnusson PK et al. The genetic architecture of economic and political preferences. *Proc Natl Acad Sci USA* 2012; **109**: 8026–8031.
16 Rietveld CA, Conley D, Eriksson N, Esko T, Medland SE, Vinkhuyzen AA et al. Replicability and robustness of genome-wide-association studies for behavioral traits. *Psychol Sci* 2014; **25**: 1975–1986.
17 Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR. Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinforma Oxf Engl* 2012; **28**: 2540–2542.
18 White KR. The relation between socioeconomic status and academic achievement. *Psychol Bull* 1982; **91**: 461–481.
19 Sirin SR. Socioeconomic status and academic achievement: a meta-analytic review of research. *Rev Educ Res* 2005; **75**: 417–453.
20 Plomin R, Deary IJ. Genetics and intelligence differences: five special findings. *Mol Psychiatry* 2014. doi:10.1038/mp.2014.105.

21 Wray NR, Lee SH, Mehta D, Vinkhuyzen AA, Dudbridge F, Middeldorp CM. Polygenic methods and their application to psychiatric disorders and related traits. *J Child Psychol Psychiatry* 2014; **55**: 1068–1087.

22 Plomin R, Simpson MA. The future of genomics for developmentalists. *Dev Psychopathol* 2013; **25**: 1263–1278.

23 Dudbridge F. Power and predictive accuracy of polygenic risk scores. *PLoS Genet* 2013; **9**: e1003348.

24 Bartels M, Rietveld MJH, Van Baal GCM, Boomsma DI. Heritability of educational achievement in 12-year-olds and the overlap with cognitive ability. *Twin Res Hum Genet* 2002; **5**: 544–553.

25 Deary IJ, Strand S, Smith P, Fernandes C. Intelligence and educational achievement. *Intelligence* 2007; **35**: 13–21.

26 Benyamin B, Pourcain B, Davis OS, Davies G, Hansell NK, Brion MJ *et al.* Childhood intelligence is heritable, highly polygenic and associated with FNBP1L. *Mol Psychiatry* 2014; **19**: 253–258.

27 Trzaskowski M, Harlaar N, Arden R, Krapohl E, Rimfeld K, McMillan A *et al.* Genetic influence on family socioeconomic status and children's intelligence. *Intelligence* 2014; **42**: 83–88.

28 Haworth CMA, Davis OSP, Plomin R. Twins Early Development Study (TEDS): a genetically sensitive investigation of cognitive and behavioral development from childhood to young adulthood. *Twin Res Hum Genet* 2013; **16**: 117–125.

29 Oliver BR, Plomin R. Twins' Early Development Study (TEDS): a multivariate, longitudinal genetic investigation of language, cognition and behavior problems from childhood through adolescence. *Twin Res Hum Genet* 2007; **10**: 96–105.

30 Kovas Y, Haworth CMA, Dale PS, Plomin R. The genetic and environmental origins of learning abilities and disabilities in the early school years. *Monogr Soc Res Child Dev* 2007; **72**, vii, 1–144.

31 Trzaskowski M, Eley TC, Davis OS, Doherty SJ, Hanscombe KB, Meaburn EL *et al.* First genome-wide association study on anxiety-related behaviours in childhood. *PLoS One* 2013; **8**: e58676.

32 Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet* 2009; **5**: e1000529.

33 Fellay J, Shianna KV, Ge D, Colombo S, Ledergerber B, Weale M *et al.* A whole-genome association study of major determinants for host control of HIV-1. *Science* 2007; **317**: 944–947.

34 Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet* 2006; **2**: e190.

35 National pupil database - GOV.UK. Available at https://www.gov.uk/government/collections/national-pupil-database.

36 Hanscombe KB, Trzaskowski M, Haworth CM, Davis OS, Dale PS, Plomin R *et al.* Socioeconomic status (SES) and children's intelligence (IQ): in a UK-representative sample SES moderates the environmental, not genetic, effect on IQ. *PLoS One* 2012; **7**: e30320.

37 Raven J, Court J, Raven J. *Manual for Raven's Progressive Matrices and Vocabulary Scales.* Oxford University Press: Oxford, 1996.

38 Raven J, Court J, Raven J. *Manual for Raven's Progressive Matrices.* HK Lewis: London, 1998.

39 Raven J, Raven J, Court J. *Mill Hill Vocabulary Scale.* OPP, 1998.

40 Office for National Statistics, United Kingdom. *Standard occupational classification 2000: Structure and description of unit groups.* Stationery Office: London, UK, 2000.

41 Office for National Statistics, United Kingdom. *Standard Occupational Classification 2000: The coding index.* Stationery Office: London, UK, 2000.

42 Trzaskowski M, Yang J, Visscher PM, Plomin R. DNA evidence for strong genetic stability and increasing heritability of intelligence from age 7 to 12. *Mol Psychiatry* 2013; **19**: 380–384.

43 R Core Team. *R: A Language and Environment for Statistical Computing.* R Foundation for Statistical Computing, 2013 Available at http://www.R-project.org/.

44 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* 2007; **81**: 559–575.

45 Plomin R, DeFries JC, Knopik VS, Neiderhiser JM. *Behavioral Genetics.* 6th ed. Worth Publishers: New York, 2013.

46 Johnson W, McGue M, Iacono WG. Genetic and environmental influences on academic achievement trajectories during adolescence. *Dev Psychol* 2006; **42**: 514–532.

47 Wray NR, Yang J, Hayes BJ, Price AL, Goddard ME, Visscher PM. Pitfalls of predicting complex traits from SNPs. *Nat Rev Genet* 2013; **14**: 507–515.

48 Krapohl E, Rimfeld K, Shakeshaft NG, Trzaskowski M, McMillan A, Pingault JB *et al.* The high heritability of educational achievement reflects many genetically influenced traits, not just intelligence. *Proc Natl Acad Sci USA* 2014; **111**: 15273–15278.

49 De Zeeuw EL, van Beijsterveldt CEM, Glasner TJ, Bartels M, Ehli EA, Davies GE *et al.* Polygenic scores associated with educational attainment in adults predict educational achievement and ADHD symptoms in children. *Am J Med Genet B Neuropsychiatr Genet* 2014; **165**: 510–520.

50 Kapur S, Phillips AG, Insel TR. Why has it taken so long for biological psychiatry to develop clinical tests and what to do about it? *Mol Psychiatry* 2012; **17**: 1174–1179.

51 Gaugler T, Klei L, Sanders SJ, Bodea CA, Goldberg AP, Lee AB *et al.* Most genetic risk for autism resides with common variation. *Nat Genet* 2014; **46**: 881–885.

52 Yang L, Neale BM, Liu L, Lee SH, Wray NR, Ji N *et al.* Polygenic transmission and complex neuro developmental network for attention deficit hyperactivity disorder: genome-wide association study of both common and rare variants. *Am J Med Genet B Neuropsychiatr Genet* 2013; **162B**: 419–430.

53 Falconer DS. *Introduction to Quantitative Genetics.* Ronald: New York, NY, USA, 1960.

54 Haworth CMA, Davis OSP. From observational to dynamic genetics. *Front Genet* 2014; **5**: 6.

55 Heath AC, Berg K, Eaves LJ, Solaas MH, Corey LA, Sundet J *et al.* Education policy and the heritability of educational attainment. *Nature* 1985; **314**: 734–736.

56 Tucker-Drob EM, Briley DA. Continuity of genetic and environmental influences on cognition across the life span: a meta-analysis of longitudinal twin and adoption studies. *Psychol Bull* 2014; **140**: 949–979.

57 Plomin R, DeFries JC, Loehlin JC. Genotype-environment interaction and correlation in the analysis of human behavior. *Psychol Bull* 1977; **84**: 309–322.

58 Kendler KS, Baker JH. Genetic influences on measures of the environment: a systematic review. *Psychol Med* 2007; **37**: 615–626.

59 Plomin R. *Genetics and Experience: The Interplay Between Nature and Nurture.* Sage Publications, Inc: Thousand Oaks, CA, USA, 1994.

60 Plomin R, Bergeman CS. The nature of nurture: genetic influence on 'environmental' measures. *Behav Brain Sci* 1991; **14**: 373–386.

61 Vinkhuyzen AAE, van der Sluis S, de Geus EJC, Boomsma DI, Posthuma D. Genetic influences on 'environmental' factors. *Genes Brain Behav* 2010; **9**: 276–287.

62 Young MD. *The Rise of the Meritocracy.* Thames and Hudson: London, UK, 1958.

63 Saunders P. *Social Mobility Delusions: Why So Much of What Politicians Say about Social Mobility in Britain Is Wrong, Misleading or Unreliable.* Civitas: London, UK, 2012.

64 Asbury K, Plomin R. *G is for Genes: The Impact of Genetics on Education and Achievement.* Wiley-Blackwell: Chichester, UK, 2013.

Supplementary Information accompanies the paper on the Molecular Psychiatry website (http://www.nature.com/mp)

# 6 Widespread covariation of early environmental exposures and trait-associated polygenic variation

Krapohl, E., Hannigan, L. J., Pingault, J.-B., Patel, H., Kadeva, N., Curtis, C. J. C., Breen, G., Newhouse, S., Eley, T. C., O'Reilly, P. F., and Plomin, R. (2017a). Widespread covariation of early environmental exposures and trait-associated polygenic variation. *Proceedings of the National Academy of Sciences*

Supplementary material: `http://www.pnas.org/content/suppl/2017/10/11/1707178114.DCSupplemental`

Article redacted from e-thesis due to journal copyright restrictions.
Please read article on journal website:
http://www.pnas.org/content/114/44/11727.full.pdf

# 7 General summary and discussion

Multi-variables approaches for trait prediction are a response to the ubiquitous poly-genicity of complex traits and the sharing of genetic effects across traits. This thesis described a series of multi-variable approaches for polygenic trait prediction and for investigating gene-environment correlation in the context of polygenic prediction models.

The following sections briefly summarise the findings, and discuss general limitations and possible future directions.

## 7.1 Summary of findings

**Chapter 3** presented a phenome-wide analysis of genome-wide polygenic scores, a systematic investigation of associations between 13 polygenic scores for cognitive, psychiatric and anthropometric traits and 50 behavioural outcomes measured in adolescence. Cognitive and educational polygenic scores yielded stronger predictions than psychiatric polygenic scores in the UK-representative sample of adolescents.

One of the main limitations of the research presented in chapter 3 is the very limited statistical power of the genetic predictors. The statistical power of the polygenic scores is mainly a function of the statistical power of the discovery GWAS (i.e. variance explained and sample size) on which they are based (Dudbridge, 2013). In the years between the two studies presented in chapters 3 and 4, larger GWAS for the same traits were published (Davies et al., 2016; Demontis et al., 2017; Okbay et al., 2016; PGC, 2017), some of which yielded increased single score prediction in the second study (chapter 4). This trajectory is likely to continue with ever-larger GWAS being published (Dudbridge, 2013; Visscher et al., 2017).

Despite the modest effect sizes of the polygenic score predictions in chapter 3, quantile analyses suggested the research potential of stratifying individuals by poly-genic score quantile. Although this study investigated multiple traits generating a profile of cross-trait polygenic associations, it did so as a series of univariate models rather than within one prediction model.

The study presented in **chapter 4** employed a multi-polygenic score (MPS) approach to increase prediction by exploiting the joint power of multiple discovery GWAS. The MPS approach improved prediction over the best single-score models for three child

outcomes, educational achievement, cognitive ability, and BMI. The findings suggest the usefulness of multi-variable approaches for trait prediction, which is projected to increase with the publication of new and larger discovery GWAS.

The prediction models presented in chapter 4 were only an initial illustration or proof-of-concept of the MPS approach. Although the MPS models showed better out-of-sample prediction than the best single-score model, it remains to be seen whether this will replicate for different outcome traits and different samples. For a more parsimonious and computationally straightforward replication, future research in other samples, could construct a simple multiple regression model using the top five predictors selected by the current analyses.

Importantly, this approach is entirely focused on prediction and does not elucidate mechanisms. This is due to two things: First, the magnitude of prediction of a given polygenic score is a function of both the genetic correlation between the predictor and the outcome trait as well as the statistical power of the discovery study the predictor is based on. Second, the observation of a genetic correlation does not in itself provide insight into underlying mechanisms or aetiology.

**Chapter 5** provided DNA-based evidence for substantial genetic influence on differences in children's educational achievement and its association with family socioeconomic status (SES). GCTA analyses showed that the estimated proportion of variance tagged by the sampled SNPs was 31% in educational achievement, and 20% in family SES. These analyses demonstrated the ability of DNA-based methods to explore the genetic architecture of extended phenotypes such as family SES that cannot be detected by traditional variance/covariance estimation methods such as the twin method that rely on known kinship relatedness. GCTA results showed that SNPs that are associated with both family SES and exam scores capture about half of their phenotypic correlation. Moreover, about one-third of this genetic association also extended to children's intelligence, but two-thirds of the genetic association between family SES and exam scores were statistically independent of intelligence.

Providing converging evidence, the polygenic score analyses showed that SNPs associated with total years of education in adulthood discovered by an independent GWA meta-analysis explained up to 3% of the variance in children's educational achievement, and up to 2% of the variance after controlling for intelligence.

The findings of a genetic correlation between parental socioeconomic status and children's educational achievement estimated by GCTA as well as the polygenic score

for years of schooling predicting children's educational achievement have since been replicated in a Dutch and in a different UK cohort (Davies et al., 2015; de Zeeuw et al., 2014).

The analyses presented in **chapter 6** investigated covariation between trait-associated polygenic variation and environmental exposures, controlling for overall genetic relatedness using a genomic-relatedness-matrix restricted/residual maximum-likelihood model.

First, the findings replicated the general finding of genotype-environment correlation for a wide range of 'environmental exposures' that classic epidemiological research often conceptualises as purely environmental in origin. The findings also show that associations between parenting behaviours and child outcomes are partially explained by genetic factors. Second, the findings show that DNA variants identified by trait GWAS are associated with parental behaviours and in part underlie correlations between parental behaviour and children's outcomes.

The analyses show that as genetic similarity in conventionally unrelated individuals increases, so does similarity of environmental exposure, likely partially explained by parents providing both environment and genotype for offspring. Existing polygenic score prediction models do not take this into account.

The findings suggest that incorporating genetic variation associated with established environmental risk or protective factors might improve genotype-based trait prediction. Although the findings provide evidence for the relevance of gene-environment correlation for polygenic trait prediction methods, they are not informative about the mechanisms involved.

## 7.2 General limitations

In addition to the limitations discussed in the papers and above, there are several general limitations that merit reiteration.

### 7.2.1 Correlation does not equal causation

Importantly no directionality or causality can be inferred from the observed genetic correlations. It is important to reiterate that genetic correlation between two phenotypes or environments does not imply that the same biological processes give rise to variation in both (i.e. 'pleiotropy').

Genetic correlation can arise from pleiotropy, the phenomenon of multiple traits being associated with the same genetic variant or genomic region (Solovieff et al., 2013; Visscher and Yang, 2016; Wright, 1984). Two types of pleiotropy are conventionally distinguished. In 'biological pleiotropy' a single genomic locus directly influences multiple traits via independent causal paths. 'Mediated pleiotropy' refers to a single genetic process, which results in a chain of events with one trait lying on another trait's causal path (Solovieff et al., 2013). Under the assumption of pure 'mediated pleiotropy', the genetic variants (or set of variants) can be used as an instrumental variable to test the causal effect of the first on the second trait (Catapano and Ference, 2015; Davey-Smith and Hemani, 2014; Hemani et al., 2016; Kathiresan, 2015).

Genetic correlation can also arise from generational effects, or what could be termed cross-generational pleiotropy. Parents pass on both genotype and environment to the offspring generation. Therefore, a genetic correlation between trait A and environment B or trait C in the offspring generation can arise from: a genetic predisposition for trait A being passed on from parents, and a predisposition for trait A causing parents to provide an environment B to the offspring or induce trait C in the offspring (Richmond et al., 2017; Zhang et al., 2015).

Genetic correlation between traits can also be the result of non-random mating or what could be termed cross-mate pleiotropy. Non-random mating where mates are correlated for two genetically influenced traits, e.g. taller individuals are more likely to mate with smarter individuals (Keller et al., 2013), induces LD between the loci associated with the assorted traits in the next generation. The offspring of cross-trait assortatively-mating parents will show an increased co-occurrence of alleles that are associated with the two traits. Effects of non-random mating will accrue over generations until reaching an equilibrium (Falconer and Mackay, 1996).

Importantly, the mechanisms described above are not mutually exclusive. For instance, the observed genetic correlation between parental characteristics and children's genotypes found in chapters 5 and 6 are likely a result of an amalgam of pleiotropic, parental, and assortative mating effects. Therefore, it is important to

keep in mind that an observation of a genetic correlation between two traits is not per se informative about the underlying mechanisms. However, as shown in this thesis, genetic correlation can be used in polygenic prediction models while remaining agnostic about underlying mechanisms; and genotype-based complex trait prediction could be used for prevention and intervention strategies long before causal mechanisms are determined.

### 7.2.2 Trait prediction limited to European-ancestry samples

The present investigations are limited to only one population, with the study sample being British and the GWAS discovery samples largely based on European samples. Although there exists suggestive evidence for some genetic risk variants to be shared across multiple ancestries (de Candia et al., 2013), it has been shown that polygenic scores created based on European GWAS are biased by genetic drift in other populations, with biases in any direction possible (Martin et al., 2017). Based on this, the predictions observed in the current investigations are unlikely to reliably replicate in samples with non-European samples.

More generally, the limited portability of trait-SNP association estimated by European populations to other populations suggests the need for and potential gain from more generalised genomic prediction methods based on the inclusion of more diverse populations to allow for prediction in populations with non-European ancestry. Methodological approaches leveraging trans-ethnic information will likely play an important role in this endeavour (Coram et al., 2017).

### 7.2.3 Additive effects of common variants

All methods employed in the current investigations are limited to detecting additive effects of common genetic variants measured on (or imputed from) conventional genotyping arrays. There is evidence that additive effects explain the majority of the total genetic variance (Hill et al., 2008; Visscher et al., 2017; Zhu et al., 2015). While non-additive effects likely exist, the power to detect these is a function of the proportion of the variance they capture, reducing the probability that they will be detected unless variants have intermediate frequency.

It remains unknown how much of the additive genetic variation is explained by low-frequency ($<1\%$) genetic variants (Visscher et al., 2017). Whole genome sequencing data in large samples will allow for explicitly estimating the contribution of low-

frequency variants to trait variation. First evidence suggest that contributions from the low frequency spectrum will likely differ by phenotype (Loh et al., 2015a; Moser et al., 2015; Ripke et al., 2013; van Rheenen et al., 2016; Visscher et al., 2017), and might deviate from expectations of an evolutionarily neutral model (Yang et al., 2015).

### 7.2.4 Upper limit of SNP-based trait prediction

The accuracy of SNP-heritability ($h^2_{SNP}$) estimates is of conceptual and logistic importance for genotype-based trait prediction, which currently typically relies on common SNPs and therefore is subject to the upper limit of total trait variation explained by common SNPs. A general methodological limitation of $h^2_{SNP}$ estimation methods is whether its prior assumptions about the distribution of heritability across the genome hold. To avoid overfitting when estimating $h^2_{SNP}$, it is necessary to make strong assumptions about marker effect sizes.

The conventionally employed GCTA and closely related methods assume a uniform Gaussian distribution of effect sizes for each marker (Bulik-Sullivan et al., 2015b; Loh et al., 2015b; Yang et al., 2011a, 2013; Zhou and Stephens, 2014). Recently, it has been shown how $h^2_{SNP}$ varies with minor allele frequency (MAF), linkage disequilibrium (LD) and genotype certainty, differently to what is assumed by GCTA (Speed et al., 2017). When effect size was modelled as a function of local LD (as well as marker quality score), $h^2_{SNP}$ estimates increased. An intuitive way of thinking about this is that if variants in a genomic region are highly correlated they are likely to all tag the same causal variant, therefore a model that gives equal weight to all markers might overestimate contributions of markers in high and thereby underestimate contributions of markers in low LD regions. It was also shown that markers with lower population frequency contribute less to heritability than assumed by the GCTA model, suggesting estimated $h^2_{SNP}$, might be increased simply by revising genotype scaling. Further improvements in $h^2_{SNP}$ estimates might be possible by incorporating functional annotations such as proximity to coding regions in the effect size prior (Speed et al., 2017).

The implication of this is that the additive effects of the typically genotyped (and imputed) common markers might be higher and the 'missing heritability' gap between $h^2$ and $h^2_{SNP}$ might be smaller than previously assumed. Higher $h^2_{SNP}$ would imply more scope for GWAS and resulting polygenic score prediction whose ceiling for the total effect size is defined by $h^2_{SNP}$.

### 7.2.5    Limits of individual-level prediction

Although genotype-based prediction has value for prediction at the individual level for Mendelian and increasingly also for oligogenic traits; this is not currently the case for most polygenic traits. Immense logistic and conceptual challenges will have to be overcome before genotype-based prediction accuracy ($R^2$) will approach trait heritability ($h^2$). These challenges include increasing discovery samples, typing rare variation, and improving parameter specification, as outlined in the introduction and the section above. What is more, even under the hypothetical scenario of $R^2=h^2$, individual-level prediction will be moderate unless $h^2$ is extremely high.

The role of stochasticity in individual differences is often understated, with the focus being placed on either genetic or environmental factors of variance shared between individuals. Phenotypic discordance of bilateral organs of the same organism or of genetically-identical twins growing up in the same environment explains substantial proportions of variation in many phenotypes (Davey-Smith, 2011; Plomin et al., 2001; Plomin and Daniels, 1987). This variance component is labelled 'non-shared environment' by twin studies. However, because of its non-stability across development and in the absence of identification of its sources, it can arguably most parsimoniously be conceptualised as random or stochastic events. This randomness induces unreliability rather than systematic biases on population-level prediction estimates. However, stochastic variability places an inherent limit on individual-level prediction.

## 7.3    Implications and possible future directions

Several implications and potential future areas of research arise from the work presented in this thesis, a selection of which is outlined in the following.

### 7.3.1    Stratified rather than personalised prediction

The substantial stochastic element in human trait variation implies that, strictly speaking, actual individual-level trait prediction is impossible, even if all genetic effects were perfectly estimated. That is, predicted outcomes will not reliably differ from one individual to the next, but between sub-groups of individuals. Nevertheless, genotype-based prediction might allow for meaningful stratification, i.e. division of individuals into groups with practically distinct risk (or resilience) profiles. Therefore, group-level prediction can achieve high individual-level relevance.

Stratification of prevention strategies could be one of the most effective intermediate-term implications of genotype-based prediction. Because discrimination between individuals within the middle of the normal distribution of complex traits or common diseases is unobtainable in the near future, identifying the tails of the polygenic distribution is a crude but pragmatic approach for increasing statistical power to predict practically relevant differences between groups of individuals.

For instance, the analyses in chapter 4 found one and a half-grade difference between individuals' educational achievement in the top versus bottom 10% group of the multi-polygenic score distribution. More generally, research has shown that polygenic scores explain a sufficient proportion of variation to stratify groups, for example, samples with the highest and lowest risk. Polygenic score models have for instance been used to stratify by disease onset (Ahn et al., 2016; Escott-Price et al., 2015; Power et al., 2017), disease comorbidity (Hamshere et al., 2013; Wiste et al., 2014), treatment resistance (Frank et al., 2015), and treatment response (Musci et al., 2016). Increased relative risk reduction through intervention in high genetic risk compared to low genetic risk groups has been shown for complex traits, for instance in cardiovascular disease (Khera et al., 2016; Mega et al., 2015; Natarajan et al., 2017).

By increasing the overall proportion of individual differences explained by the prediction model, multi-polygenic score approaches have the potential to minimise the sub-groups between which differences in outcome can be reliably predicted.

### 7.3.2 Prediction beyond family history

If the goal is (early) prediction, the ultimate test of the value of genomic prediction must be whether it does better than other readily available information. Family history, which reflects genetic and non-genetic influences, has so far been superior to genomic data in predicting most complex traits and diseases, except for heritable traits with very low prevalence (Chatterjee et al., 2013; Cornelis et al., 2015; Do et al., 2012; McGrath et al., 2013).

However, prediction based on family risk is family-general, whereas genotype-based trait prediction is individual-specific. Therefore, polygenic prediction might add value to family history, especially in cases where familial risk is high. For instance, familial studies of breast cancer risk suggest the utility of SNP-based prediction in the context of elevated familial risk (Li et al., 2017; Mavaddat et al., 2013; Muranen

et al., 2016). In BRCA1/2-negative women with a family history of breast cancer, polygenic risk scores yielded some incremental risk prediction beyond family history, suggesting potential for more effective or more tailored prevention strategies.

Another example of individual trait-associated genomic variation explaining within-family phenotypic variation comes from the realm of educational attainment. Specifically, studies have shown that individuals' polygenic score for education predicted phenotypic deviation from parental or sibling educational attainment (Ayorech et al., 2017; Domingue et al., 2015).

More generally, with only half of the alleles transmitted from each parent, the highly polygenic nature of many traits suggests that parent and offspring genetic propensities may differ considerably. Moreover, family history might not always be straightforwardly available or unobtainable for instance for phenotypes such as drug response. What is more, differential genetic liability of family members such as siblings might be especially relevant when family history indicates elevated family-level risk.

### 7.3.3 Leveraging within-family genetic variation for prediction

Of all genetic variation in the population, 50% occurs between siblings within families. Although this fact has been extensively used by heritability estimation methods, arguably, it has not been fully exploited for genomic prediction. For instance, genotypes of dizygotic twin represent random samples from the same pool of parental genotypes. Therefore, any variance in outcome traits predicted by genetic differences between dizygotic twins would be highly valuable as it would be incremental to variance predicted by family-level factors.

To achieve this, a mixed model could be fitted according to this general schema: $Y_{ij} = \mu + G_{ij}^{trait} + F_i + E_{ij}$, with $\mu$ being the average phenotypic value, the fixed effect $G_{ij}^{trait}$ measuring individual-level effects (deviation of the $j^{th}$ individual's trait $Y$ from the average for the $i^{th}$ family) of a polygenic score, $F_i$ being the family-specific random effect (difference between the average trait at family $i$ and the average trait in the sample), and $E_{ij}$ containing residual variance.

This allows for testing whether genetic differences between dizygotic twins predict outcome beyond (genetic and non-genetic) family-level factors. Formally, a likelihood-ratio-test could test the null hypothesis that G is zero, by comparing the likelihoods of models GFE and FE. The model could be extended to include further

parameters: $A_{ij}^{individual}$, capturing individual-level additive genetic variance that is not captured by the polygenic score; and $A_{ij}^{family}$ measuring family-level additive genetic variance that is not captured by the polygenic predictor.

A similar model has been used before: The precursor for $h_{SNP}^2$ estimation in unrelated individuals, which used genome-wide markers in sibling pairs to estimate heritability (Visscher et al., 2006). Although this pedigree-based $h_{SNP}^2$ method can estimate heritability contributed by different parts of the genome (genome, chromosome, gene etc.), its purpose is not trait prediction. In contrast, in the above model $G^{trait}$ contains externally discovered trait-associated variation in form of a polygenic score, estimating individual-level trait-specific prediction incremental to family-level factors. This approach could also be expanded to include multiple polygenic scores.

### 7.3.4  Expansion of the multi-polygenic score approach

The multi-polygenic score approach presented in chapter 4 could be applied to a wide range of outcomes and different types of samples, including disease classification in case-control samples. It could be further extended by considering the findings from chapter 6 that individuals' trait-associated polygenic variation captures variance in established environmental risk and protective factors.

Using cardiovascular disease (CAD) as an example, the following illustrates how a multi-polygenic score approach combined with a consideration of genetic influences on environmental risk factors might be used to increase genetic risk prediction.

A survey of $\sim$56,000 found that a polygenic risk score of 50 externally discovered risk variants predicted incident of CAD, with people in the top 20% of the genetic risk distribution having 91% higher relative risk compared to those in the bottom 20% (Khera et al., 2016). The study also showed that in any genetic risk quintile, 'adherence to healthy lifestyle' was associated with relative risk reductions in CAD event rates, with up to 50% in participants within the highest genetic risk score quintile. The healthy life style factors included 'no obesity' (BMI<30), 'no current smoking', weakly physical activity, and a healthy diet, with the first two being the strongest risk predictors. These factors have been shown to be under genetic influence (Bauman et al., 2012; Locke et al., 2015; Maes et al., 1997; Rose et al., 2009; Tobacco and Genetics Consortium, 2010; Zaitlen et al., 2013). Therefore, incorporating known genetic variation associated with such factors into the prediction model

in form of multiple polygenic scores may be of interest for early prediction and possibly for stratifying intervention strategies as a function of genetic predispositions for 'environmental' or 'lifestyle' risk factors.

This kind of approach might be fruitful for range of outcomes that have genetic and environmental risk factors and where typically only a single genetic predictor is included in genetic risk prediction models.

## 7.4 General conclusion

The polygenic and pleiotropic architecture of complex traits and their correlation with environmental exposures greatly complicates trait prediction. Joint modelling of many phenotypes (including environmental 'phenotypes') and many genetic variants offers many advantages over the traditional approach of considering only marginal single-variant single-phenotype associations. Multi-variable approaches have the potential to provide more complete individual inventories of genetic propensity for traits and environmental exposures.

Taken together, the current investigations illustrate the value of multi-variable approaches to complex trait prediction and investigation of genotype-environment correlation, as well as their current limitations and future promise.

Ultimately, whole genome sequencing (Levy et al., 2007) could allow for predicting trait variation from all sequence variants. In the interim, genome variation can be approximated by genotyped (and imputed) common SNPs. Benefits from predicting complex traits from polygenic models can be gained long before causal mechanisms are identified. Genomic profiles available at birth might facilitate early prediction including specifying individual risk within familial risk.

Genetic prediction could allow for cost-effective prevention and intervention strategies by targeting subsets of the population for whom relative risk reduction is highest. Combining genomic with more classical predictors such as family history and environmental risk factors rather than substituting them, is likely to allow for more fine-grained prediction. In face of the complex genetic architecture and considerable stochasticity of trait variation, the challenge lays in minimising the population strata between which trait variation can be reliably predicted. Multi-variable approaches represent a pragmatic tool for this task.

# References

Afifi, T. O., Mota, N. P., Dasiewicz, P., MacMillan, H. L., and Sareen, J. (2012). Physical Punishment and Mental Disorders: Results From a Nationally Representative US Sample. *Pediatrics*, 130(2):184–192.

Ahn, K., An, S. S., Shugart, Y. Y., and Rapoport, J. L. (2016). Common polygenic variation and risk for childhood-onset schizophrenia. *Molecular Psychiatry*, 21(1):94–96.

Ainsworth, J. W. (2002). Why Does It Take a Village? The Mediation of Neighborhood Effects on Educational Achievement. *Social Forces*, 81(1):117–152.

Avinun, R. and Knafo, A. (2013). Parenting as a Reaction Evoked by Children's Genotype A Meta-Analysis of Children-as-Twins Studies. *Personality and Social Psychology Review*, page 1088868313498308.

Ayorech, Z., Krapohl, E., Plomin, R., and von Stumm, S. (2017). Genetic Influence on Intergenerational Educational Attainment. *Psychological Science*, page 0956797617707270.

Barton, N. H. (1990). Pleiotropic models of quantitative variation. *Genetics*, 124(3):773–782.

Baselmans, B. M., Jansen, R., Dongen, J. v., Bao, Y., Smart, M., Kumari, M., Abdellaoui, A., Weijer, M. P. v. d., Willemsen, G., Hottenga, J. J., Consortium, B., Consortium, S. S. G. A., Geus, E. J. d., Boomsma, D. I., Nivard, M. G., and Bartels, M. (2017). Multivariate Genome-Wide and Integrated Transcriptome and Epigenome-Wide Analyses of the Well-being Spectrum. *bioRxiv*, page 115915.

Bauman, A. E., Reis, R. S., Sallis, J. F., Wells, J. C., Loos, R. J., and Martin, B. W. (2012). Correlates of physical activity: why are some people physically active and others not? *The Lancet*, 380(9838):258–271.

Bender, H. L., Allen, J. P., McElhaney, K. B., Antonishak, J., Moore, C. M., O&amp, H., apos, Kelly, B., and Davis, S. M. (2007). Use of harsh physical discipline and developmental outcomes in adolescence. *Development and Psychopathology*, 19(1):227–242.

Benjamin, D. J., Cesarini, D., Loos, M. J. H. M. v. d., Dawes, C. T., Koellinger, P. D., Magnusson, P. K. E., Chabris, C. F., Conley, D., Laibson, D., Johannesson, M., and Visscher, P. M. (2012). The genetic architecture of economic and political preferences. *Proceedings of the National Academy of Sciences*, 109(21):8026–8031.

Bhattacharjee, S., Rajaraman, P., Jacobs, K. B., Wheeler, W. A., Melin, B. S., Hartge, P., GliomaScan Consortium, Yeager, M., Chung, C. C., Chanock, S. J., and Chatterjee, N. (2012). A subset-based approach improves power and interpretation for the combined analysis of genetic association studies of heterogeneous traits. *American Journal of Human Genetics*, 90(5):821–835.

Bolormaa, S., Pryce, J. E., Reverter, A., Zhang, Y., Barendse, W., Kemper, K., Tier, B., Savin, K., Hayes, B. J., and Goddard, M. E. (2014). A Multi-Trait, Meta-analysis for Detecting Pleiotropic Polymorphisms for Stature, Fatness and Reproduction in Beef Cattle. *PLOS Genetics*, 10(3):e1004198.

Bulik-Sullivan, B., Finucane, H. K., Anttila, V., Gusev, A., Day, F. R., Consortium, R., Consortium, P. G., 3, G. C. f. A. N. o. t. W. T. C., Perry, J. R. B., Patterson, N., Robinson, E., Daly, M. J., Price, A. L., and Neale, B. M. (2015a). An Atlas of Genetic Correlations across Human Diseases and Traits. *bioRxiv*, page 014498.

Bulik-Sullivan, B. K., Loh, P.-R., Finucane, H. K., Ripke, S., Yang, J., Schizophrenia Working Group of the Psychiatric Genomics Consortium, Patterson, N., Daly, M. J., Price, A. L., and Neale, B. M. (2015b). LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature Genetics*, 47(3):291–295.

Butcher, L. M. and Plomin, R. (2008). The Nature of Nurture: A Genomewide Association Scan for Family Chaos. *Behavior Genetics*, 38(4):361–371.

Byrne, M., Agerbo, E., Ewald, H., Eaton, W. W., and Mortensen, P. B. (2003). Parental age and risk of schizophrenia: a case-control study. *Archives of general psychiatry*, 60(7):673–678.

Calvin, C. M., Deary, I. J., Webbink, D., Smith, P., Fernandes, C., Lee, S. H., Luciano, M., and Visscher, P. M. (2012). Multivariate genetic analyses of cognition and academic achievement from two population samples of 174,000 and 166,000 school children. *Behavior Genetics*, 42(5):699–710.

Caspi, A., Houts, R. M., Belsky, D. W., Harrington, H., Hogan, S., Ramrakha, S., Poulton, R., and Moffitt, T. E. (2016). Childhood forecasting of a small segment of the population with large economic burden. *Nature Human Behaviour*, 1:0005.

Catapano, A. L. and Ference, B. A. (2015). IMPROVE-IT and genetics reaffirm the causal role of LDL in Cardiovascular Disease. *Atherosclerosis*, 241(2):498–501.

Chatterjee, N., Wheeler, B., Sampson, J., Hartge, P., Chanock, S. J., and Park, J.-H. (2013). Projecting the performance of risk prediction based on polygenic analyses of genome-wide association studies. *Nature Genetics*, 45(4):400–405.

Colen, C. G. and Ramey, D. M. (2014). Is breast truly best? Estimating the effects of breastfeeding on long-term child health and wellbeing in the United States using sibling comparisons. *Social Science & Medicine*, 109:55–65.

Coram, M. A., Fang, H., Candille, S. I., Assimes, T. L., and Tang, H. (2017). Leveraging Multi-Ethnic Evidence for Risk Assessment of Quantitative Traits in Minority Populations. *The American Journal of Human Genetics*, 0(0).

Cornelis, M. C., Zaitlen, N., Hu, F. B., Kraft, P., and Price, A. L. (2015). Genetic and environmental components of family history in type 2 diabetes. *Human Genetics*, 134(2):259–267.

Daetwyler, H. D., Villanueva, B., and Woolliams, J. A. (2008). Accuracy of predicting the genetic risk of disease using a genome-wide approach. *PloS One*, 3(10):e3395.

Danner, F. W. (2008). A National Longitudinal Study of the Association Between Hours of TV Viewing and the Trajectory of BMI Growth Among US Children. *Journal of Pediatric Psychology*, 33(10):1100–1107.

Davey-Smith, G. (2011). Epidemiology, epigenetics and the 'Gloomy Prospect': embracing randomness in population health research and practice. *International Journal of Epidemiology*, 40(3):537–562.

Davey-Smith, G. and Hemani, G. (2014). Mendelian randomization: genetic anchors for causal inference in epidemiological studies. *Human Molecular Genetics*, 23(R1):R89–R98.

Davies, G., Marioni, R. E., Liewald, D. C., Hill, W. D., Hagenaars, S. P., Harris, S. E., Ritchie, S. J., Luciano, M., Fawns-Ritchie, C., Lyall, D., Cullen, B., Cox, S. R., Hayward, C., Porteous, D. J., Evans, J., McIntosh, A. M., Gallacher, J., Craddock, N., Pell, J. P., Smith, D. J., Gale, C. R., and Deary, I. J. (2016). Genome-wide association study of cognitive functions and educational attainment in UK Biobank (N=112,Äâ151). *Molecular Psychiatry*, 21(6):758–767.

Davies, N. M., Hemani, G., Timpson, N. J., Windmeijer, F., and Davey Smith, G. (2015). The role of common genetic variation in educational attainment and income: evidence from the National Child Development Study. *Scientific Reports*, 5.

Davis, O. S. P., Band, G., Pirinen, M., Haworth, C. M. A., Meaburn, E. L., Kovas, Y., Harlaar, N., Docherty, S. J., Hanscombe, K. B., Trzaskowski, M., Curtis, C. J. C., Strange, A., Freeman, C., Bellenguez, C., Su, Z., Pearson, R., Vukcevic, D., Langford, C., Deloukas, P., Hunt, S., Gray, E., Dronov, S., Potter, S. C., Tashakkori-Ghanbaria, A., Edkins, S., Bumpstead, S. J., Blackwell, J. M., Bramon, E., Brown, M. A., Casas, J. P., Corvin, A., Duncanson, A., Jankowski, J. A. Z., Markus, H. S., Mathew, C. G., Palmer, C. N. A., Rautanen, A., Sawcer, S. J., Trembath, R. C., Viswanathan, A. C., Wood, N. W., Barroso, I., Peltonen, L., Dale, P. S., Petrill, S. A., Schalkwyk, L. S., Craig, I. W., Lewis, C. M., Price, T. S., The Wellcome Trust Case Control Consortium, Donnelly, P., Plomin, R., and Spencer, C. C. A. (2014). The correlation between reading and mathematics ability at age twelve has a substantial genetic component. *Nature Communications*, 5.

de Candia, T. R., Lee, S. H., Yang, J., Browning, B. L., Gejman, P. V., Levinson, D. F., Mowry, B. J., Hewitt, J. K., Goddard, M. E., O'Donovan, M. C., Purcell, S. M., Posthuma, D., Visscher, P. M., Wray, N. R., and Keller, M. C. (2013). Additive Genetic Variation in Schizophrenia Risk Is Shared by Populations of African and European Descent. *The American Journal of Human Genetics*, 93(3):463–470.

de Kluiver, H., Buizer-Voskamp, J., Dolan, C., and Boomsma, D. (2017). Paternal age and psychiatric disorders: A review. *American journal of medical genetics. Part B, Neuropsychiatric genetics: the official publication of the International Society of Psychiatric Genetics*, 174(3):202.

de Zeeuw, E. L., van Beijsterveldt, C. E., Glasner, T. J., Bartels, M., Ehli, E. A., Davies, G. E., Hudziak, J. J., Social Science Genetic Association Consortium, Rietveld, C. A., Groen-Blokhuis, M. M., Hottenga, J. J., de Geus, E. J., and Boomsma, D. I. (2014). Polygenic scores associated with educational attainment in adults predict educational achievement and ADHD symptoms in children. *American Journal of Medical Genetics Part B: Neuropsychiatric Genetics*, 165(6):510–520.

Demontis, D., Walters, R. K., Martin, J., Mattheisen, M., Als, T. D., Agerbo, E., Belliveau, R., Bybjerg-Grauholm, J., Bækved-Hansen, M., Cerrato, F., Chambert, K., Churchhouse, C., Dumont, A., Eriksson, N., Gandal, M., Goldstein, J., Grove, J., Hansen, C. S., Hauberg, M., Hollegaard, M., Howrigan, D. P., Huang, H., Maller, J., Martin, A. R., Moran, J., Pallesen, J., Palmer, D. S., Pedersen, C. B., Pedersen, M. G., Poterba, T., Poulsen, J. B., Ripke, S., Robinson, E. B., Satterstrom, F. K., Stevens, C., Turley, P., Won, H., Con, A. W. G. o. t. P. G., Lifecourse &amp, E., Epidemiology (EAGLE), G., Team, a. R., Andreassen, O. A., Burton, C., Boomsma, D., Cormand, B., Dalsgaard, S., Franke, B., Gelernter, J., Geschwind, D., Hakonarson, H., Haavik, J., Kranzler, H., Kuntsi, J., Langley, K., Lesch, K.-P., Middeldorp, C., Reif, A., Rohde, L. A., Roussos, P., Schachar, R., Sklar, P., Sonuga-Barke, E., Sullivan, P. F., Thapar, A., Tung, J., Waldman, I., Nordentoft, M., Hougaard, D. M., Werge, T., Mors, O., Mortensen, P. B., Daly, M. J., Faraone, S. V., Børglum, A. D., and Neale, B. M. (2017). Discovery Of The First Genome-Wide Significant Risk Loci For ADHD. *bioRxiv*, page 145581.

Do, C. B., Hinds, D. A., Francke, U., and Eriksson, N. (2012). Comparison of Family History and SNPs for Predicting Risk of Complex Disease. *PLoS Genetics*, 8(10).

Domingue, B. W., Belsky, D. W., Conley, D., Harris, K. M., and 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):2332858415599972.

D'Onofrio, B. M., Hulle, C. A. V., Waldman, I. D., Rodgers, J. L., Rathouz, P. J., and Lahey, B. B. (2007). Causal Inferences Regarding Prenatal Alcohol Exposure and Childhood Externalizing Problems. *Archives of General Psychiatry*, 64(11):1296–1304.

D'Onofrio, B. M., Rickert, M. E., Frans, E., Kuja-Halkola, R., Almqvist, C., Sjölander, A., Larsson, H., and Lichtenstein, P. (2014). Paternal Age at Childbearing and Offspring Psychiatric and Academic Morbidity. *JAMA Psychiatry*, 71(4):432–438.

D'Onofrio, B. M., Singh, A. L., Iliadou, A., Lambe, M., Hultman, C. M., Grann, M., Neiderhiser, J. M., and Lichtenstein, P. (2010a). Familial Confounding of the Association Between Maternal Smoking During Pregnancy and Offspring Criminality: A Population-Based Study in Sweden. *Archives of General Psychiatry*, 67(5):529–538.

D'Onofrio, B. M., Singh, A. L., Iliadou, A., Lambe, M., Hultman, C. M., Neiderhiser, J. M., Långström, N., and Lichtenstein, P. (2010b). A Quasi-Experimental Study of Maternal Smoking During Pregnancy and Offspring Academic Achievement. *Child Development*, 81(1):80–100.

Dudbridge, F. (2013). Power and predictive accuracy of polygenic risk scores. *PLoS Genet*, 9(3):e1003348.

Duncan, L., Yilmaz, Z., Gaspar, H., Walters, R., Goldstein, J., Anttila, V., Bulik-Sullivan, B. K., Ripke, S., Thornton, L., Hinney, A., Daly, M., Sullivan, P. F., Zeggini, E., Breen, G., Bulik, C. M., Duncan, L., Yilmaz, Z., Gaspar, H., Walters, R., Goldstein, J., Anttila, V., Bulik-Sullivan, B., Ripke, S., Adan, R., Alfredsson, L., Ando, T., Andreassen, O., Aschauer, H., Baker, J., Barrett, J., Bencko, V., Bergen, A., Berrettini, W., Birgegård, A., Boni, C., Perica, V. B., Brandt, H., Burghardt, R., Carlberg, L., Cassina, M., Cesta, C., Cichon, S., Clementi, M., Cohen-Woods, S., Coleman, J., Cone, R., Courtet, P., Crawford, S., Crow, S., Crowley, J., Danner, U., Davis, O., de Zwaan, M., Dedoussis, G., Degortes, D., DeSocio, J., Dick, D., Dikeos, D., Dina, C., Ding, B., Dmitrzak-Weglarz, M., Docampo, E., Egberts, K., Ehrlich, S., Escaramís, G., Esko, T., Espeseth, T., Estivill, X., Favaro, A., Fernández-Aranda, F., Fichter, M., Finan, C., Fischer, K., Floyd, J., Föcker, M., Foretova, L., Forzan, M., Fox, C., Franklin, C., Gaborieau, V., Gallinger, S., Gambaro, G., Giegling, I., Gonidakis, F., Gorwood, P., Gratacos, M., Guillaume, S., Guo, Y., Hakonarson, H., Halmi, K., Harrison, R., Hatzikotoulas, K., Hauser, J., Hebebrand, J., Helder, S., Hendriks, J., Herms, S., Herpertz-Dahlmann, B., Herzog, W., Hilliard, C., Huckins, L., Hudson, J., Huemer, J., Imgart, H., Inoko, H., Jall, S., Jamain, S., Janout, V., Jiménez-Murcia, S., Johnson, C., Jordan, J., Julià, A., Juréus, A., Kalsi, G., Kaplan, A., Kaprio, J., Karhunen, L., Karwautz, A., Kas, M., Kaye, W., Kennedy, M., Kennedy, J., Keski-Rahkonen, A., Kiezebrink, K., Kim, Y.-R., Klareskog, L., Klump, K., Knudsen, G. P., Koeleman, B., Koubek, D., La Via, M., Landén, M., Le Hellard, S., Leboyer, M., Levitan, R., Li, D., Lichtenstein, P., Lilenfeld, L., Lissowska, J., Lundervold, A., Magistretti, P., Maj, M., Mannik, K., Marsal, S., Kaminska, D., Martin, N., Mattingsdal, M., McDevitt, S., McGuffin, P., Merl, E., Metspalu, A., Meulenbelt, I., Micali, N., Mitchell, J., Mitchell, K., Monteleone, P., Monteleone, A. M., Montgomery, G., Mortensen, P., Munn-Chernoff, M., Müller, T., Nacmias, B., Navratilova, M., Nilsson, I., Norring, C., Ntalla, I., Ophoff, R., O'Toole, J., Palotie, A., Pantel, J., Papezova, H., Parker, R., Pinto, D., Rabionet, R., Raevuori, A., Rajewski, A., Ramoz, N., Rayner, N. W., Reichborn-Kjennerud, T., Ricca, V., Ripatti, S., Ritschel, F., Roberts, M., Rotondo, A., Rujescu, D., Rybakowski, F., Santonastaso, P., Scherag, A., Scherer, S., Schmidt, U., Schork, N., Schosser, A., Scott, L., Seitz, J., Slachtova, L., Sladek, R., Slagboom, P. E., 't Landt, M. S.-O., Slopien, A., Smith, T., Soranzo, N., Sorbi, S., Southam, L., Steen, V., Strengman, E., Strober, M., Szatkiewicz, J., Szeszenia-Dabrowska, N., Tachmazidou, I., Tenconi, E., Tortorella, A., Tozzi, F., Treasure, J., Tschöp, M., Tsitsika, A., Tziouvas, K., van Elburg, A., van Furth, E., Wade, T., Wagner, G., Walton, E., Watson, H., Wichmann, H.-E., Widen, E., Woodside, D. B., Yanovski, J., Yao, S.,

Zerwas, S., Zipfel, S., Thornton, L., Hinney, A., Daly, M., Sullivan, P. F., Zeggini, E., Breen, G., and Bulik, C. M. (2017). Significant Locus and Metabolic Genetic Correlations Revealed in Genome-Wide Association Study of Anorexia Nervosa. *American Journal of Psychiatry*, page appi.ajp.2017.16121402.

Eamon, M. K. (2005). Social-Demographic, School, Neighborhood, and Parenting Influences on the Academic Achievement of Latino Young Adolescents. *Journal of Youth and Adolescence*, 34(2):163–174.

Eichler, E. E., Flint, J., Gibson, G., Kong, A., Leal, S. M., Moore, J. H., and Nadeau, J. H. (2010). Missing heritability and strategies for finding the underlying causes of complex disease. *Nature Reviews Genetics*, 11(6):446–450.

Escott-Price, V., International Parkinson's Disease Genomics Consortium, Nalls, M. A., Morris, H. R., Lubbe, S., Brice, A., Gasser, T., Heutink, P., Wood, N. W., Hardy, J., Singleton, A. B., Williams, N. M., and IPDGC consortium members (2015). Polygenic risk of Parkinson disease is correlated with disease age at onset. *Annals of Neurology*, 77(4):582–591.

Euesden, J., Lewis, C. M., and O'Reilly, P. F. (2015). PRSice: Polygenic Risk Score software. *Bioinformatics*, 31(9):1466–1468.

Evans, D. M., Visscher, P. M., and Wray, N. R. (2009). Harnessing the information contained within genome-wide association studies to improve individual prediction of complex disease risk. *Human Molecular Genetics*, 18(18):3525–3531.

Evenhouse, E. and Reilly, S. (2005). Improved Estimates of the Benefits of Breastfeeding Using Sibling Comparisons to Reduce Selection Bias. *Health Services Research*, 40(6p1):1781–1802.

Falconer, D. and Mackay, T. (1996). *Introduction to Quantitative Genetics.* Vol 4. Longman: Harlow, UK.

Ferreira, M. A. R. and Purcell, S. M. (2009). A multivariate test of association. *Bioinformatics*, 25(1):132–133.

Fisher, R. (1918). The Correlation between Relatives on the Supposition of Mendelian Inheritance. *Philosophical Transactions of the Royal Society of Edinburgh*, 52:399–433.

Frank, J., Lang, M., Witt, S. H., Strohmaier, J., Rujescu, D., Cichon, S., Degenhardt, F., Nöthen, M. M., Collier, D. A., Ripke, S., Naber, D., and Rietschel, M. (2015). Identification of increased genetic risk scores for schizophrenia in treatment-resistant patients. *Molecular Psychiatry*, 20(2):150–151.

Garner, C. L. and Raudenbush, S. W. (1991). Neighborhood Effects on Educational Attainment: A Multilevel Analysis. *Sociology of Education*, 64(4):251–262.

Garrod, A. E. (1902). The incidence of alkaptonuria: A study in chemical individuality. *The Lancet*, 160(4137):1616–1620.

Ge, T., Chen, C.-Y., Neale, B. M., Sabuncu, M. R., and Smoller, J. W. (2017). Phenome-wide heritability analysis of the UK Biobank. *PLOS Genetics*, 13(4):e1006711.

Gentile, D. A., Lynch, P. J., Linder, J. R., and Walsh, D. A. (2004). The effects of violent video game habits on adolescent hostility, aggressive behaviors, and school performance. *Journal of Adolescence*, 27(1):5–22.

Gershoff, E. T. (2002). Corporal punishment by parents and associated child behaviors and experiences: A meta-analytic and theoretical review. *Psychological Bulletin*, 128(4):539–579.

Giambartolomei, C., Vukcevic, D., Schadt, E. E., Franke, L., Hingorani, A. D., Wallace, C., and Plagnol, V. (2014). Bayesian Test for Colocalisation between Pairs of Genetic Association Studies Using Summary Statistics. *PLOS Genetics*, 10(5):e1004383.

Goddard, M. (2009). Genomic selection: prediction of accuracy and maximisation of long term response. *Genetica*, 136(2):245–257.

Hamshere, M. L., Langley, K., Martin, J., Agha, S. S., Stergiakouli, E., Anney, R. J. L., Buitelaar, J., Faraone, S. V., Lesch, K. P., Neale, B. M., Franke, B., Barke, E., Asherson, P., Merwood, A., Kuntsi, J., Medland, S. E., Ripke, S., Steinhausen, H. C., Freitag, C., Reif, A., Renner, T. J., Romanos, M., Romanos, J., Warnke, A., Meyer, J., Palmason, H., Vasquez, A. A., Lambregts-Rommelse, N., Roeyers, H., Biederman, J., Doyle, A. E., Hakonarson, H., Rothenberger, A., Banaschewski, T., Oades, R. D., McGough, J. J., Kent, L., Williams, N., Owen, M. J., Holmans, P., O'Donovan, M. C., and Thapar, A. (2013). High loading of polygenic risk for ADHD in children with comorbid aggression. *AMERICAN JOURNAL OF PSYCHIATRY*, 170(8):909–916.

Harden, K. P., Lynch, S. K., Turkheimer, E., Emery, R. E., D'Onofrio, B. M., Slutske, W. S., Waldron, M. D., Heath, A. C., Statham, D. J., and Martin, N. G. (2007). A Behavior Genetic Investigation of Adolescent Motherhood and Offspring Mental Health Problems. *Journal of abnormal psychology*, 116(4):667–683.

Hemani, G., Zheng, J., Wade, K. H., Laurin, C., Elsworth, B., Burgess, S., Bowden, J., Langdon, R., Tan, V., Yarmolinsky, J., Shihab, H. A., Timpson, N., Evans, D. M., Relton, C., Martin, R. M., Smith, G. D., Gaunt, T. R., Haycock, P. C., and Collaboration, T. M.-B. (2016). MR-Base: a platform for systematic causal inference across the phenome using billions of genetic associations. *bioRxiv*, page 078972.

Henderson, C. R., Kempthorne, O., Searle, S. R., and Krosigk, C. M. v. (1959). The Estimation of Environmental and Genetic Trends from Records Subject to Culling. *Biometrics*, 15(2):192–218.

Hill, W. D., Hagenaars, S. P., Marioni, R. E., Harris, S. E., Liewald, D. C. M., Davies, G., Okbay, A., McIntosh, A. M., Gale, C. R., and Deary, I. J. (2016). Molecular Genetic Contributions to Social Deprivation and Household Income in UK Biobank. *Current Biology*, 0(0).

Hill, W. G., Goddard, M. E., and Visscher, P. M. (2008). Data and Theory Point to Mainly Additive Genetic Variance for Complex Traits. *PLOS Genetics*, 4(2):e1000008.

Huizink, A. C. and Mulder, E. J. H. (2006). Maternal smoking, drinking or cannabis use during pregnancy and neurobehavioral and cognitive functioning in human offspring. *Neuroscience & Biobehavioral Reviews*, 30(1):24–41.

Jago, R., Baranowski, T., Baranowski, J. C., Thompson, D., and Greaves, K. A. (2005). BMI from 3–6,Ääy of age is predicted by TV viewing and physical activity, not diet. *International Journal of Obesity*, 29(6):557–564.

Janecka, M., Haworth, C. M. A., Ronald, A., Krapohl, E., Happé, F., Mill, J., Schalkwyk, L. C., Fernandes, C., Reichenberg, A., and Rijsdijk, F. (2017). Paternal Age Alters Social Development in Offspring. *Journal of the American Academy of Child & Adolescent Psychiatry*, 56(5):383–390.

Johnson, T. (2013). gtx: Genetics ToolboX.

Kathiresan, S. (2015). Developing Medicines That Mimic the Natural Successes of the Human Genome. *Journal of the American College of Cardiology*, 65(15):1562–1566.

Keller, M. C., Garver-Apgar, C. E., Wright, M. J., Martin, N. G., Corley, R. P., Stallings, M. C., Hewitt, J. K., and Zietsch, B. P. (2013). The Genetic Correlation between Height and IQ: Shared Genes or Assortative Mating? *PLoS Genet*, 9(4):e1003451.

Kendler, K. S. and Baker, J. H. (2007). Genetic influences on measures of the environment: a systematic review. *Psychological Medicine*, 37(05):615–626.

Kendler, K. S., Gardner, C. O., Annas, P., Neale, M. C., Eaves, L. J., and Lichtenstein, P. (2008). A longitudinal twin study of fears from middle childhood to early adulthood: evidence for a developmentally dynamic genome. *Archives of General Psychiatry*, 65(4):421–429.

Khera, A. V., Emdin, C. A., Drake, I., Natarajan, P., Bick, A. G., Cook, N. R., Chasman, D. I., Baber, U., Mehran, R., Rader, D. J., Fuster, V., Boerwinkle, E., Melander, O., Orho-Melander, M., Ridker, P. M., and Kathiresan, S. (2016). Genetic Risk, Adherence to a Healthy Lifestyle, and Coronary Disease. *New England Journal of Medicine*, 375(24):2349–2358.

Klahr, A. M. and Burt, S. A. (2014). Elucidating the etiology of individual differences in parenting: A meta-analysis of behavioral genetic research. *Psychological Bulletin*, 140(2):544–586.

Knopik, V. S., Heath, A. C., Jacob, T., Slutske, W. S., Bucholz, K. K., Madden, P. a. F., Waldron, M., and Martin, N. G. (2006). Maternal alcohol use disorder and offspring ADHD: disentangling genetic and environmental effects using a children-of-twins design. *Psychological Medicine*, 36(10):1461–1471.

Knox, M. (2010). On Hitting Children: A Review of Corporal Punishment in the United States. *Journal of Pediatric Health Care*, 24(2):103–107.

Korte, A., Vilhjálmsson, B. J., Segura, V., Platt, A., Long, Q., and Nordborg, M. (2012). A mixed-model approach for genome-wide association studies of correlated traits in structured populations. *Nature genetics*, 44(9):1066–1071.

Kovas, Y. and Plomin, R. (2006). Generalist genes: implications for the cognitive sciences. *Trends in Cognitive Sciences*, 10(5):198–203.

Krapohl, E., Euesden, J., Zabaneh, D., Pingault, J.-B., Rimfeld, K., von Stumm, S., Dale, P. S., Breen, G., O'Reilly, P. F., and Plomin, R. (2016). Phenome-wide analysis of genome-wide polygenic scores. *Molecular Psychiatry*, 21(9):1188–1193.

Krapohl, E., Hannigan, L. J., Pingault, J.-B., Patel, H., Kadeva, N., Curtis, C. J. C., Breen, G., Newhouse, S., Eley, T. C., O'Reilly, P. F., and Plomin, R. (2017a). Widespread covariation of early environmental exposures and trait-associated polygenic variation. *Proceedings of the National Academy of Sciences*.

Krapohl, E., Patel, H., Newhouse, S., Curtis, C., von Stumm, S., Dale, P., Zabaneh, D., Breen, G., O'Reilly, P., and Plomin, R. (2017b). Multi-polygenic score approach to trait prediction. *Molecular psychiatry*.

Krapohl, E. and Plomin, R. (2016). Genetic link between family socioeconomic status and children's educational achievement estimated from genome-wide SNPs. *Molecular Psychiatry*, 21(3):437–443.

Krapohl, E., Rimfeld, K., Shakeshaft, N. G., Trzaskowski, M., McMillan, A., Pingault, J.-B., Asbury, K., Harlaar, N., Kovas, Y., Dale, P. S., and Plomin, R. (2014). The high heritability of educational achievement reflects many genetically influenced traits, not just intelligence. *Proceedings of the National Academy of Sciences*, 111(42):15273–15278.

Larsson, H., Sariaslan, A., Långström, N., D'Onofrio, B., and Lichtenstein, P. (2014). Family income in early childhood and subsequent attention deficit/hyperactivity disorder: a quasi-experimental study. *Journal of Child Psychology and Psychiatry, and Allied Disciplines*, 55(5):428–435.

Lee, S. H. and van der Werf, J. H. (2006). An efficient variance component approach implementing an average information REML suitable for combined LD and linkage mapping with a general complex pedigree. *Genetics Selection Evolution*, 38:25.

Lee, S. H., Yang, J., Goddard, M. E., Visscher, P. M., and Wray, N. R. (2012). Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinformatics*, 28(19):2540–2542.

Leventhal, T. and Brooks-Gunn, J. (2000). The neighborhoods they live in: The effects of neighborhood residence on child and adolescent outcomes. *Psychological Bulletin*, 126(2):309–337.

Levy, S., Sutton, G., Ng, P. C., Feuk, L., Halpern, A. L., Walenz, B. P., Axelrod, N., Huang, J., Kirkness, E. F., Denisov, G., Lin, Y., MacDonald, J. R., Pang, A. W. C., Shago, M., Stockwell, T. B., Tsiamouri, A., Bafna, V., Bansal, V., Kravitz, S. A., Busam, D. A., Beeson, K. Y., McIntosh, T. C., Remington, K. A., Abril, J. F., Gill, J., Borman, J., Rogers, Y.-H., Frazier, M. E., Scherer, S. W., Strausberg, R. L., and Venter, J. C. (2007). The diploid genome sequence of an individual human. *PLoS biology*, 5(10):e254.

Li, H., Feng, B., Miron, A., Chen, X., Beesley, J., Bimeh, E., Barrowdale, D., John, E. M., Daly, M. B., Andrulis, I. L., Buys, S. S., Kraft, P., kConFab investigators, Thorne, H., Chenevix-Trench, G., Southey, M. C., Antoniou, A. C., James, P. A., Terry, M. B., Phillips, K.-A., Hopper, J. L., Mitchell, G., and Goldgar, D. E. (2017). Breast cancer risk prediction using a polygenic risk score in the familial setting: a prospective study from the Breast Cancer Family Registry and kConFab. *Genetics in Medicine: Official Journal of the American College of Medical Genetics*, 19(1):30–35.

Lichtenstein, P., Yip, B. H., Björk, C., Pawitan, Y., Cannon, T. D., Sullivan, P. F., and Hultman, C. M. (2009). Common genetic determinants of schizophrenia and bipolar disorder in Swedish families: a population-based study. *The Lancet*, 373(9659):234–239.

Liu, C., Dupuis, J., Larson, M. G., Cupples, L. A., Ordovas, J. M., Vasan, R. S., Meigs, J. B., Jacques, P. F., and Levy, D. (2015). Revisiting heritability accounting for shared environmental effects and maternal inheritance. *Human Genetics*, 134(2):169–179.

Locke, A. E., Kahali, B., Berndt, S. I., Justice, A. E., Pers, T. H., Day, F. R., Powell, C., Vedantam, S., Buchkovich, M. L., Yang, J., Croteau-Chonka, D. C., Esko, T., Fall, T., Ferreira, T., Gustafsson, S., Kutalik, Z., Luan, J., Mägi, R., Randall, J. C., Winkler, T. W., Wood, A. R., Workalemahu, T., Faul, J. D., Smith, J. A., Hua Zhao, J., Zhao, W., Chen, J., Fehrmann, R., Hedman, Å. K., Karjalainen, J., Schmidt, E. M., Absher, D., Amin, N., Anderson, D., Beekman, M., Bolton, J. L., Bragg-Gresham, J. L., Buyske, S., Demirkan, A., Deng, G., Ehret, G. B., Feenstra, B., Feitosa, M. F., Fischer, K., Goel, A., Gong, J., Jackson, A. U., Kanoni, S., Kleber, M. E., Kristiansson, K., Lim, U., Lotay, V., Mangino, M., Mateo Leach, I., Medina-Gomez, C., Medland, S. E., Nalls, M. A., Palmer, C. D., Pasko, D., Pechlivanis, S., Peters, M. J., Prokopenko, I., Shungin, D., Stančáková, A., Strawbridge, R. J., Ju Sung, Y., Tanaka, T., Teumer, A.,

Trompet, S., van der Laan, S. W., van Setten, J., Van Vliet-Ostaptchouk, J. V., Wang, Z., Yengo, L., Zhang, W., Isaacs, A., Albrecht, E., Ärnlöv, J., Arscott, G. M., Attwood, A. P., Bandinelli, S., Barrett, A., Bas, I. N., Bellis, C., Bennett, A. J., Berne, C., Blagieva, R., Blüher, M., Böhringer, S., Bonnycastle, L. L., Böttcher, Y., Boyd, H. A., Bruinenberg, M., Caspersen, I. H., Ida Chen, Y.-D., Clarke, R., Warwick Daw, E., de Craen, A. J. M., Delgado, G., Dimitriou, M., Doney, A. S. F., Eklund, N., Estrada, K., Eury, E., Folkersen, L., Fraser, R. M., Garcia, M. E., Geller, F., Giedraitis, V., Gigante, B., Go, A. S., Golay, A., Goodall, A. H., Gordon, S. D., Gorski, M., Grabe, H.-J., Grallert, H., Grammer, T. B., Gräßler, J., Grönberg, H., Groves, C. J., Gusto, G., Haessler, J., Hall, P., Haller, T., Hallmans, G., Hartman, C. A., Hassinen, M., Hayward, C., Heard-Costa, N. L., Helmer, Q., Hengstenberg, C., Holmen, O., Hottenga, J.-J., James, A. L., Jeff, J. M., Johansson, Å., Jolley, J., Juliusdottir, T., Kinnunen, L., Koenig, W., Koskenvuo, M., Kratzer, W., Laitinen, J., Lamina, C., Leander, K., Lee, N. R., Lichtner, P., Lind, L., Lindström, J., Sin Lo, K., Lobbens, S., Lorbeer, R., Lu, Y., Mach, F., Magnusson, P. K. E., Mahajan, A., McArdle, W. L., McLachlan, S., Menni, C., Merger, S., Mihailov, E., Milani, L., Moayyeri, A., Monda, K. L., Morken, M. A., Mulas, A., Müller, G., Müller-Nurasyid, M., Musk, A. W., Nagaraja, R., Nöthen, M. M., Nolte, I. M., Pilz, S., Rayner, N. W., Renstrom, F., Rettig, R., Ried, J. S., Ripke, S., Robertson, N. R., Rose, L. M., Sanna, S., Scharnagl, H., Scholtens, S., Schumacher, F. R., Scott, W. R., Seufferlein, T., Shi, J., Vernon Smith, A., Smolonska, J., Stanton, A. V., Steinthorsdottir, V., Stirrups, K., Stringham, H. M., Sundström, J., Swertz, M. A., Swift, A. J., Syvänen, A.-C., Tan, S.-T., Tayo, B. O., Thorand, B., Thorleifsson, G., Tyrer, J. P., Uh, H.-W., Vandenput, L., Verhulst, F. C., Vermeulen, S. H., Verweij, N., Vonk, J. M., Waite, L. L., Warren, H. R., Waterworth, D., Weedon, M. N., Wilkens, L. R., Willenborg, C., Wilsgaard, T., Wojczynski, M. K., Wong, A., Wright, A. F., Zhang, Q., The LifeLines Cohort Study, Brennan, E. P., Choi, M., Dastani, Z., Drong, A. W., Eriksson, P., Franco-Cereceda, A., Gådin, J. R., Gharavi, A. G., Goddard, M. E., Handsaker, R. E., Huang, J., Karpe, F., Kathiresan, S., Keildson, S., Kiryluk, K., Kubo, M., Lee, J.-Y., Liang, L., Lifton, R. P., Ma, B., McCarroll, S. A., McKnight, A. J., Min, J. L., Moffatt, M. F., Montgomery, G. W., Murabito, J. M., Nicholson, G., Nyholt, D. R., Okada, Y., Perry, J. R. B., Dorajoo, R., Reinmaa, E., Salem, R. M., Sandholm, N., Scott, R. A., Stolk, L., Takahashi, A., Tanaka, T., Hooft, F. M. v., Vinkhuyzen, A. A. E., Westra, H.-J., Zheng, W., Zondervan, K. T., The ADIPOGen Consortium, The AGEN-BMI Working Group, The CARDIOGRAMplusC4D Consortium, The CKDGen Consortium, The Glgc, The Icbp, The MAGIC Investigators, The MuTHER Consortium, The MIGen Consortium, The PAGE Consortium, The ReproGen Consortium, The GENIE Consortium, The International Endogene Consortium, Heath, A. C., Arveiler, D., Bakker, S. J. L., Beilby, J., Bergman, R. N., Blangero, J., Bovet, P., Campbell, H., Caulfield, M. J., Cesana, G., Chakravarti, A., Chasman, D. I., Chines, P. S., Collins, F. S., Crawford, D. C., Adrienne Cupples, L., Cusi, D., Danesh, J., de Faire, U., den Ruijter, H. M., Dominiczak, A. F., Erbel, R., Erdmann, J., Eriksson, J. G., Farrall, M., Felix, S. B., Ferrannini, E., Ferrières, J., Ford, I., Forouhi, N. G., Forrester, T., Franco, O. H., Gansevoort, R. T., Gejman, P. V., Gieger, C., Gottesman, O., Gudnason, V., Gyllensten, U., Hall, A. S., Harris, T. B., Hattersley, A. T., Hicks, A. A., Hindorff, L. A., Hingorani, A. D., Hofman, A., Homuth, G., Kees Hovingh, G., Humphries, S. E., Hunt, S. C., Hyppönen, E., Illig, T., Jacobs, K. B., Jarvelin, M.-R., Jöckel, K.-H., Johansen, B., Jousilahti, P., Wouter Jukema, J., Jula, A. M., Kaprio, J., Kastelein, J. J. P., Keinanen-Kiukaanniemi, S. M., Kiemeney, L. A., Knekt, P., Kooner, J. S., Kooperberg, C., Kovacs, P., Kraja, A. T., Kumari, M., Kuusisto, J., Lakka, T. A., Langenberg, C., Le Marchand, L., Lehtimäki, T., Lyssenko, V., Männistö, S., Marette, A., Matise, T. C., McKenzie, C. A., McKnight, B., Moll, F. L., Morris, A. D., Morris, A. P., Murray, J. C., Nelis, M., Ohlsson, C., Oldehinkel, A. J., Ong, K. K., Madden, P. A. F., Pasterkamp, G., Peden, J. F., Peters, A., Postma, D. S., Pramstaller, P. P., Price, J. F., Qi, L., Raitakari, O. T., Rankinen, T., Rao, D. C., Rice, T. K., Ridker, P. M., Rioux, J. D., Ritchie, M. D., Rudan, I., Salomaa, V., Samani, N. J., Saramies, J., Sarzynski, M. A., Schunkert, H.,

Schwarz, P. E. H., Sever, P., Shuldiner, A. R., Sinisalo, J., Stolk, R. P., Strauch, K., Tönjes, A., Trégouët, D.-A., Tremblay, A., Tremoli, E., Virtamo, J., Vohl, M.-C., Völker, U., Waeber, G., Willemsen, G., Witteman, J. C., Carola Zillikens, M., Adair, L. S., Amouyel, P., Asselbergs, F. W., Assimes, T. L., Bochud, M., Boehm, B. O., Boerwinkle, E., Bornstein, S. R., Bottinger, E. P., Bouchard, C., Cauchi, S., Chambers, J. C., Chanock, S. J., Cooper, R. S., de Bakker, P. I. W., Dedoussis, G., Ferrucci, L., Franks, P. W., Froguel, P., Groop, L. C., Haiman, C. A., Hamsten, A., Hui, J., Hunter, D. J., Hveem, K., Kaplan, R. C., Kivimaki, M., Kuh, D., Laakso, M., Liu, Y., Martin, N. G., März, W., Melbye, M., Metspalu, A., Moebus, S., Munroe, P. B., Njølstad, I., Oostra, B. A., Palmer, C. N. A., Pedersen, N. L., Perola, M., Pérusse, L., Peters, U., Power, C., Quertermous, T., Rauramaa, R., Rivadeneira, F., Saaristo, T. E., Saleheen, D., Sattar, N., Schadt, E. E., Schlessinger, D., Eline Slagboom, P., Snieder, H., Spector, T. D., Thorsteinsdottir, U., Stumvoll, M., Tuomilehto, J., Uitterlinden, A. G., Uusitupa, M., van der Harst, P., Walker, M., Wallaschofski, H., Wareham, N. J., Watkins, H., Weir, D. R., Wichmann, H.-E., Wilson, J. F., Zanen, P., Borecki, I. B., Deloukas, P., Fox, C. S., Heid, I. M., O'Connell, J. R., Strachan, D. P., Stefansson, K., van Duijn, C. M., Abecasis, G. R., Franke, L., Frayling, T. M., McCarthy, M. I., Visscher, P. M., Scherag, A., Willer, C. J., Boehnke, M., Mohlke, K. L., Lindgren, C. M., Beckmann, J. S., Barroso, I., North, K. E., Ingelsson, E., Hirschhorn, J. N., Loos, R. J. F., and Speliotes, E. K. (2015). Genetic studies of body mass index yield new insights for obesity biology. *Nature*, 518(7538):197–206.

Loh, P.-R., Bhatia, G., Gusev, A., Finucane, H. K., Bulik-Sullivan, B. K., Pollack, S. J., Schizophrenia Working Group of the Psychiatric Genomics Consortium, de Candia, T. R., Lee, S. H., Wray, N. R., Kendler, K. S., O'Donovan, M. C., Neale, B. M., Patterson, N., and Price, A. L. (2015a). Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nature Genetics*, 47(12):1385–1392.

Loh, P.-R., Tucker, G., Bulik-Sullivan, B. K., Vilhjálmsson, B. J., Finucane, H. K., Salem, R. M., Chasman, D. I., Ridker, P. M., Neale, B. M., Berger, B., Patterson, N., and Price, A. L. (2015b). Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nature Genetics*, 47(3):284–290.

Luciano, M., Svinti, V., Campbell, A., Marioni, R. E., Hayward, C., Wright, A. F., Taylor, M. S., Porteous, D. J., Thomson, P., Prendergast, J. G., Hastie, N. D., Farrington, S. M., Scotland, G., Dunlop, M. G., and Deary, I. J. (2015). Exome Sequencing to Detect Rare Variants Associated With General Cognitive Ability: A Pilot Study. *Twin Research and Human Genetics*, 18(02):117–125.

Lynch, S. K., Turkheimer, E., D'Onofrio, B. M., Mendle, J., Emery, R. E., Slutske, W. S., and Martin, N. G. (2006). A Genetically Informed Study of the Association Between Harsh Punishment and Offspring Behavioral Problems. *Journal of family psychology : JFP : journal of the Division of Family Psychology of the American Psychological Association (Division 43)*, 20(2):190–198.

Maes, H. H., Neale, M. C., and Eaves, L. J. (1997). Genetic and environmental factors in relative body weight and human adiposity. *Behavior Genetics*, 27(4):325–351.

Maier, R., Moser, G., Chen, G.-B., Ripke, S., Cross-Disorder Working Group of the Psychiatric Genomics Consortium, Coryell, W., Potash, J. B., Scheftner, W. A., Shi, J., Weissman, M. M., Hultman, C. M., Landén, M., Levinson, D. F., Kendler, K. S., Smoller, J. W., Wray, N. R., Lee, S. H., and Cross-Disorder Working Group of the Psychiatric Genomics Consortium (2015). Joint analysis of psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major depressive disorder. *American Journal of Human Genetics*, 96(2):283–294.

Majumdar, A., Haldar, T., Bhattacharya, S., and Witte, J. (2017). An efficient Bayesian meta-analysis approach for studying cross-phenotype genetic associations. *bioRxiv*, page 101543.

Malaspina, D. (2001). Paternal factors and schizophrenia risk: de novo mutations and imprinting. *Schizophrenia Bulletin*, 27(3):379–393.

Manolio, T. A., Collins, F. S., Cox, N. J., Goldstein, D. B., Hindorff, L. A., Hunter, D. J., Mc-Carthy, M. I., Ramos, E. M., Cardon, L. R., Chakravarti, A., Cho, J. H., Guttmacher, A. E., Kong, A., Kruglyak, L., Mardis, E., Rotimi, C. N., Slatkin, M., Valle, D., Whittemore, A. S., Boehnke, M., Clark, A. G., Eichler, E. E., Gibson, G., Haines, J. L., Mackay, T. F. C., McCarroll, S. A., and Visscher, P. M. (2009). Finding the missing heritability of complex diseases. *Nature*, 461(7265):747–753.

Marioni, R. E., Davies, G., Hayward, C., Liewald, D., Kerr, S. M., Campbell, A., Luciano, M., Smith, B. H., Padmanabhan, S., Hocking, L. J., Hastie, N. D., Wright, A. F., Porteous, D. J., Visscher, P. M., and Deary, I. J. (2014). Molecular genetic contributions to socioeconomic status and intelligence. *Intelligence*, 44:26–32.

Marioni, R. E., Ritchie, S. J., Joshi, P. K., Hagenaars, S. P., Okbay, A., Fischer, K., Adams, M. J., Hill, W. D., Davies, G., Consortium, S. S. G. A., Nagy, R., Amador, C., Läll, K., Metspalu, A., Liewald, D. C., Campbell, A., Wilson, J. F., Hayward, C., Esko, T., Porteous, D. J., Gale, C. R., Deary, I. J., Okbay, A., Beauchamp, J. P., Fontana, M. A., Lee, J. J., Pers, T. H., Rietveld, C. A., Turley, P., Chen181, G.-B., Emilsson, V., Meddens, S. F. W., Oskarsson, S., Pickrell, J. K., Thom, K., Timshel, P., Vlaming, R. d., Abdellaoui, A., Ahluwalia, T. S., Bacelis, J., Baumbach, C., Bjornsdottir, G., Brandsma, J. H., Concas, M. P., Derringer, J., Furlotte, N. A., Galesloot, T. E., Girotto, G., Gupta, R., Hall, L. M., Harris, S. E., Hofer, E., Horikoshi, M., Huffman, J. E., Kaasik, K., Kalafati, I. P., Karlsson, R., Kong, A., Lahti, J., Lee, S. J. v. d., Leeuw, C. d., Lind, P. A., Lindgren, K.-O., Liu, T., Mangino, M., Marten, J., Mihailov, E., Miller, M. B., Most, P. J. v. d., Oldmeadow, C., Payton, A., Pervjakova, N., Peyrot, W. J., Qian, Y., Raitakari, O., Rueedi, R., Salvi, E., Schmidt, B., Schraut, K. E., Shi, J., Smith, A. V., Poot, R. A., Pourcain, B. S., Teumer, A., Thorleifsson, G., Verweij, N., Vuckovic, D., Wellmann, J., Westra, H.-J., Yang, J., Zhao, W., Zhu, Z., Alizadeh, B. Z., Amin, N., Bakshi, A., Baumeister, S. E., Biino, G., Bønnelykke, K., Boyle, P. A., Campbell, H., Cappuccio, F. P., Davies, G., Neve, J.-E. D., Deloukas, P., Demuth, I., Ding, J., Eibich, P., Eisele, L., Eklund, N., Evans, D. M., Faul, J. D., Feitosa, M. F., Forstner, A. J., Gandin, I., Gunnarsson, B., Halldórsson, B. V., Harris, T. B., Heath, A. C., Hocking, L. J., Holliday, E. G., Homuth, G., Horan, M. A., Hottenga, J.-J., Jager, P. L. d., Joshi, P. K., Jugessur, A., Kaakinen, M. A., Kähönen, M., Kanoni, S., Keltigangas-Järvinen, L., Kiemeney, L. A. L. M., Kolcic, I., Koskinen, S., Kraja, A. T., Kroh, M., Kutalik, Z., Latvala, A., Launer, L. J., Lebreton, M. P., Levinson, D. F., Lichtenstein, P., Lichtner, P., Liewald, D. C. M., Study, L. C., Loukola, A., Madden, P. A., Mägi, R., Mäki-Opas, T., Marioni, R. E., Marques-Vidal, P., Meddens, G. A., McMahon, G., Meisinger, C., Meitinger, T., Milaneschi, Y., Milani, L., Montgomery, G. W., Myhre, R., Nelson, C. P., Nyholt, D. R., Ollier, W. E. R., Palotie, A., Paternoster, L., Pedersen, N. L., Petrovic, K. E., Porteous, D. J., Räikkönen, K., Ring, S. M., Robino, A., Rostapshova, O., Rudan, I., Rustichini, A., Salomaa, V., Sanders, A. R., Sarin, A.-P., Schmidt, H., Scott, R. J., Smith, B. H., Smith, J. A., Staessen, J. A., Steinhagen-Thiessen, E., Strauch, K., Terracciano, A., Tobin, M. D., Ulivi, S., Vaccargiu, S., Quaye, L., Rooij, F. J. A. v., Venturini, C., Vinkhuyzen, A. A. E., Völker, U., Völzke, H., Vonk, J. M., Vozzi, D., Waage, J., Ware, E. B., Willemsen, G., Attia, J. R., Bennett, D. A., Berger, K., Bertram, L., Bisgaard, H., Boomsma, D. I., Borecki, I. B., Bultmann, U., Chabris, C. F., Cucca, F., Cusi, D., Deary, I. J., Dedoussis, G. V., Duijn, C. M. v., Eriksson, J. G., Franke, B., Franke, L., Gasparini, P., Gejman, P. V., Gieger, C., Grabe, H.-J., Gratten, J., Groenen, P. J. F., Gudnason, V., Harst, P. v. d., Hayward, C., Hinds, D. A., Hoffmann, W., Hypponen,

E., Iacono, W. G., Jacobsson, B., Järvelin, M.-R., Jöckel, K.-H., Kaprio, J., Kardia, S. L. R., Lehtimäki, T., Lehrer, S. F., Magnusson, P. K. E., Martin, N. G., McGue, M., Metspalu, A., Pendleton, N., Penninx, B. W. J. H., Perola, M., Pirastu, N., Pirastu, M., Polasek, O., Posthuma, D., Power, C., Province, M. A., Samani, N. J., Schlessinger, D., Schmidt, R., Sørensen, T. I. A., Spector, T. D., Stefansson, K., Thorsteinsdottir, U., Thurik, A. R., Timpson, N. J., Tiemeier, H., Tung, J. Y., Uitterlinden, A. G., Vitart, V., Vollenweider, P., Weir, D. R., Wilson, J. F., Wright, A. F., Conley, D. C., Krueger, R. F., Smith, G. D., Hofman, A., Laibson, D. I., Medland, S. E., Meyer, M. N., Yang, J., Johannesson, M., Visscher, P. M., Esko, T., Koellinger, P. D., Cesarini, D., and Benjamin, D. J. (2016). Genetic variants linked to education predict longevity. *Proceedings of the National Academy of Sciences*, 113(47):13366–13371.

Martin, A. R., Gignoux, C. R., Walters, R. K., Wojcik, G. L., Neale, B. M., Gravel, S., Daly, M. J., Bustamante, C. D., and Kenny, E. E. (2017). Human Demographic History Impacts Genetic Risk Prediction across Diverse Populations. *The American Journal of Human Genetics*, 100(4):635–649.

Mavaddat, N., Pharoah, P., Hall, P., Easton, D., and Garcia-Closas, M. (2013). Abstract 2567: Genetic risk stratification for breast cancer based on a polygenic risk score and family history. *Cancer Research*, 73(8 Supplement):2567–2567.

McGrath, J. J., Mortensen, P. B., Visscher, P. M., and Wray, N. R. (2013). Where GWAS and Epidemiology Meet: Opportunities for the Simultaneous Study of Genetic and Environmental Risk Factors in Schizophrenia. *Schizophrenia Bulletin*, 39(5):955–959.

Mega, J. L., Stitziel, N. O., Smith, J. G., Chasman, D. I., Caulfield, M. J., Devlin, J. J., Nordio, F., Hyde, C. L., Cannon, C. P., Sacks, F. M., Poulter, N. R., Sever, P. S., Ridker, P. M., Braunwald, E., Melander, O., Kathiresan, S., and Sabatine, M. S. (2015). Genetic risk, coronary heart disease events, and the clinical benefit of statin therapy: an analysis of primary and secondary prevention trials. *The Lancet*, 385(9984):2264–2271.

Mendel, G. (1866). Versuche über Pflanzenhybriden. *Verhandlungen des naturforschenden Vereines in Brunn 4: 3*, 44.

Meuwissen, T. H. E., Hayes, B. J., and Goddard, M. E. (2001). Prediction of Total Genetic Value Using Genome-Wide Dense Marker Maps. *Genetics*, 157(4):1819–1829.

Moser, G., Lee, S. H., Hayes, B. J., Goddard, M. E., Wray, N. R., and Visscher, P. M. (2015). Simultaneous Discovery, Estimation and Prediction Analysis of Complex Traits Using a Bayesian Mixture Model. *PLoS Genet*, 11(4):e1004969.

Mousseau, T. A., Ritland, K., and Heath, D. D. (1998). A novel method for estimating heritability using molecular markers. *Heredity*, 80(2):218–224.

Muñoz, M., Pong-Wong, R., Canela-Xandri, O., Rawlik, K., Haley, C. S., and Tenesa, A. (2016). Evaluating the contribution of genetics and familial shared environment to common disease using the UK Biobank. *Nature Genetics*, 48(9):980–983.

Muranen, T. A., Mavaddat, N., Khan, S., Fagerholm, R., Pelttari, L., Lee, A., Aittomäki, K., Blomqvist, C., Easton, D. F., and Nevanlinna, H. (2016). Polygenic risk score is associated with increased disease risk in 52 Finnish breast cancer families. *Breast Cancer Research and Treatment*, 158(3):463–469.

Musci, R. J., Fairman, B., Masyn, K. E., Uhl, G., Maher, B., Sisto, D. Y., Kellam, S. G., and Ialongo, N. S. (2016). Polygenic Score x Intervention Moderation: an Application of Discrete-Time Survival Analysis to Model the Timing of First Marijuana Use Among Urban Youth. *Prevention Science*, pages 1–9.

Narusyte, J., Neiderhiser, J. M., D'Onofrio, B. M., Reiss, D., Spotts, E. L., Ganiban, J., and Lichtenstein, P. (2008). Testing different types of genotype-environment correlation: An extended children-of-twins model. *Developmental Psychology*, 44(6):1591–1603.

Natarajan, P., Young, R., Stitziel, N. O., Padmanabhan, S., Baber, U., Mehran, R., Sartori, S., Fuster, V., Reilly, D. F., Butterworth, A. S., Rader, D. J., Ford, I., Sattar, N., and Kathiresan, S. (2017). Polygenic Risk Score Identifies Subgroup with Higher Burden of Atherosclerosis and Greater Relative Benefit from Statin Therapy in the Primary Prevention Setting. *Circulation*, page CIRCULATIONAHA.116.024436.

Okbay, A., Beauchamp, J. P., Fontana, M. A., Lee, J. J., Pers, T. H., Rietveld, C. A., Turley, P., Chen, G.-B., Emilsson, V., Meddens, S. F. W., Oskarsson, S., Pickrell, J. K., Thom, K., Timshel, P., de Vlaming, R., Abdellaoui, A., Ahluwalia, T. S., Bacelis, J., Baumbach, C., Bjornsdottir, G., Brandsma, J. H., Pina Concas, M., Derringer, J., Furlotte, N. A., Galesloot, T. E., Girotto, G., Gupta, R., Hall, L. M., Harris, S. E., Hofer, E., Horikoshi, M., Huffman, J. E., Kaasik, K., Kalafati, I. P., Karlsson, R., Kong, A., Lahti, J., Lee, S. J. v. d., deLeeuw, C., Lind, P. A., Lindgren, K.-O., Liu, T., Mangino, M., Marten, J., Mihailov, E., Miller, M. B., van der Most, P. J., Oldmeadow, C., Payton, A., Pervjakova, N., Peyrot, W. J., Qian, Y., Raitakari, O., Rueedi, R., Salvi, E., Schmidt, B., Schraut, K. E., Shi, J., Smith, A. V., Poot, R. A., St Pourcain, B., Teumer, A., Thorleifsson, G., Verweij, N., Vuckovic, D., Wellmann, J., Westra, H.-J., Yang, J., Zhao, W., Zhu, Z., Alizadeh, B. Z., Amin, N., Bakshi, A., Baumeister, S. E., Biino, G., Bønnelykke, K., Boyle, P. A., Campbell, H., Cappuccio, F. P., Davies, G., De Neve, J.-E., Deloukas, P., Demuth, I., Ding, J., Eibich, P., Eisele, L., Eklund, N., Evans, D. M., Faul, J. D., Feitosa, M. F., Forstner, A. J., Gandin, I., Gunnarsson, B., Halldórsson, B. V., Harris, T. B., Heath, A. C., Hocking, L. J., Holliday, E. G., Homuth, G., Horan, M. A., Hottenga, J.-J., de Jager, P. L., Joshi, P. K., Jugessur, A., Kaakinen, M. A., Kähönen, M., Kanoni, S., Keltigangas-Järvinen, L., Kiemeney, L. A. L. M., Kolcic, I., Koskinen, S., Kraja, A. T., Kroh, M., Kutalik, Z., Latvala, A., Launer, L. J., Lebreton, M. P., Levinson, D. F., Lichtenstein, P., Lichtner, P., Liewald, D. C. M., Cohort Study, L., Loukola, A., Madden, P. A., Mägi, R., Mäki-Opas, T., Marioni, R. E., Marques-Vidal, P., Meddens, G. A., McMahon, G., Meisinger, C., Meitinger, T., Milaneschi, Y., Milani, L., Montgomery, G. W., Myhre, R., Nelson, C. P., Nyholt, D. R., Ollier, W. E. R., Palotie, A., Paternoster, L., Pedersen, N. L., Petrovic, K. E., Porteous, D. J., Räikkönen, K., Ring, S. M., Robino, A., Rostapshova, O., Rudan, I., Rustichini, A., Salomaa, V., Sanders, A. R., Sarin, A.-P., Schmidt, H., Scott, R. J., Smith, B. H., Smith, J. A., Staessen, J. A., Steinhagen-Thiessen, E., Strauch, K., Terracciano, A., Tobin, M. D., Ulivi, S., Vaccargiu, S., Quaye, L., van Rooij, F. J. A., Venturini, C., Vinkhuyzen, A. A. E., Völker, U., Völzke, H., Vonk, J. M., Vozzi, D., Waage, J., Ware, E. B., Willemsen, G., Attia, J. R., Bennett, D. A., Berger, K., Bertram, L., Bisgaard, H., Boomsma, D. I., Borecki, I. B., Bültmann, U., Chabris, C. F., Cucca, F., Cusi, D., Deary, I. J., Dedoussis, G. V., van Duijn, C. M., Eriksson, J. G., Franke, B., Franke, L., Gasparini, P., Gejman, P. V., Gieger, C., Grabe, H.-J., Gratten, J., Groenen, P. J. F., Gudnason, V., van der Harst, P., Hayward, C., Hinds, D. A., Hoffmann, W., Hyppönen, E., Iacono, W. G., Jacobsson, B., Järvelin, M.-R., Jöckel, K.-H., Kaprio, J., Kardia, S. L. R., Lehtimäki, T., Lehrer, S. F., Magnusson, P. K. E., Martin, N. G., McGue, M., Metspalu, A., Pendleton, N., Penninx, B. W. J. H., Perola, M., Pirastu, N., Pirastu, M., Polasek, O., Posthuma, D., Power, C., Province, M. A., Samani, N. J., Schlessinger, D., Schmidt, R., Sørensen, T. I. A., Spector, T. D., Stefansson, K., Thorsteinsdottir, U., Thurik, A. R., Timpson, N. J., Tiemeier, H., Tung, J. Y., Uitterlinden, A. G., Vitart, V., Vollenweider, P., Weir, D. R.,

Wilson, J. F., Wright, A. F., Conley, D. C., Krueger, R. F., Davey Smith, G., Hofman, A., Laibson, D. I., Medland, S. E., Meyer, M. N., Yang, J., Johannesson, M., Visscher, P. M., Esko, T., Koellinger, P. D., Cesarini, D., and Benjamin, D. J. (2016). Genome-wide association study identifies 74 loci associated with educational attainment. *Nature*, 533(7604):539–542.

Palla, L. and Dudbridge, F. (2015). A Fast Method that Uses Polygenic Scores to Estimate the Variance Explained by Genome-wide Marker Panels and the Proportion of Variants Affecting a Trait. *American Journal of Human Genetics*, 97(2):250–259.

PGC (2013). Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *The Lancet*, 381(9875):1371–1379.

PGC (2017). Meta-analysis of GWAS of over 16,000 individuals with autism spectrum disorder highlights a novel locus at 10q24.32 and a significant overlap with schizophrenia. *Molecular Autism*, 8(1):21.

Pickrell, J. K., Berisa, T., Liu, J. Z., Ségurel, L., Tung, J. Y., and Hinds, D. A. (2016). Detection and interpretation of shared genetic influences on 42 human traits. *Nature Genetics*, 48(7):709–717.

Plomin, R., Asbury, K., and Dunn, J. (2001). Why are Children in the Same Family So Different? Nonshared Environment a Decade Later. *The Canadian Journal of Psychiatry*, 46(3):225–233.

Plomin, R. and Bergeman, C. S. (1991). The nature of nurture: Genetic influence on "environmental" measures. *Behavioral and Brain Sciences*, 14(03):373–386.

Plomin, R. and Daniels, D. (1987). Why are children in the same family so different from one another? *Behavioral and Brain Sciences*, 10(1):1–16.

Plomin, R. and DeFries, J. C. (1979). Multivariate behavioral genetic analysis of twin data on scholastic abilities. *Behavior Genetics*, 9(6):505–517.

Polderman, T. J. C., Benyamin, B., de Leeuw, C. A., Sullivan, P. F., van Bochoven, A., Visscher, P. M., and Posthuma, D. (2015). Meta-analysis of the heritability of human traits based on fifty years of twin studies. *Nature Genetics*, advance online publication.

Power, R. A., Steinberg, S., Bjornsdottir, G., Rietveld, C. A., Abdellaoui, A., Nivard, M. M., Johannesson, M., Galesloot, T. E., Hottenga, J. J., Willemsen, G., Cesarini, D., Benjamin, D. J., Magnusson, P. K. E., Ullén, F., Tiemeier, H., Hofman, A., van Rooij, F. J. A., Walters, G. B., Sigurdsson, E., Thorgeirsson, T. E., Ingason, A., Helgason, A., Kong, A., Kiemeney, L. A., Koellinger, P., Boomsma, D. I., Gudbjartsson, D., Stefansson, H., and Stefansson, K. (2015). Polygenic risk scores for schizophrenia and bipolar disorder predict creativity. *Nature Neuroscience*, 18(7):953–955.

Power, R. A., Tansey, K. E., Buttenschøn, H. N., Cohen-Woods, S., Bigdeli, T., Hall, L. S., Kutalik, Z., Lee, S. H., Ripke, S., Steinberg, S., Teumer, A., Viktorin, A., Wray, N. R., Arolt, V., Baune, B. T., Boomsma, D. I., Børglum, A. D., Byrne, E. M., Castelao, E., Craddock, N., Craig, I. W., Dannlowski, U., Deary, I. J., Degenhardt, F., Forstner, A. J., Gordon, S. D., Grabe, H. J., Grove, J., Hamilton, S. P., Hayward, C., Heath, A. C., Hocking, L. J., Homuth, G., Hottenga, J. J., Kloiber, S., Krogh, J., Landén, M., Lang, M., Levinson, D. F., Lichtenstein, P., Lucae, S., MacIntyre, D. J., Madden, P., Magnusson, P. K. E., Martin, N. G., McIntosh, A. M., Middeldorp, C. M., Milaneschi, Y., Montgomery, G. W., Mors, O., Müller-Myhsok, B., Nyholt, D. R., Oskarsson, H., Owen, M. J., Padmanabhan, S., Penninx, B. W. J. H., Pergadia, M. L., Porteous, D. J., Potash, J. B., Preisig, M., Rivera, M., Shi, J., Shyn, S. I., Sigurdsson, E., Smit, J. H., Smith, B. H., Stefansson, H., Stefansson, K., Strohmaier, J., Sullivan, P. F., Thomson, P., Thorgeirsson,

T. E., Van der Auwera, S., Weissman, M. M., CONVERGE Consortium, CARDIoGRAM Consortium, GERAD1 Consortium, Breen, G., and Lewis, C. M. (2017). Genome-wide Association for Major Depression Through Age at Onset Stratification: Major Depressive Disorder Working Group of the Psychiatric Genomics Consortium. *Biological Psychiatry*, 81(4):325–335.

Power, R. A., Wingenbach, T., Cohen-Woods, S., Uher, R., Ng, M. Y., Butler, A. W., Ising, M., Craddock, N., Owen, M. J., Korszun, A., Jones, L., Jones, I., Gill, M., Rice, J. P., Maier, W., Zobel, A., Mors, O., Placentino, A., Rietschel, M., Lucae, S., Holsboer, F., Binder, E. B., Keers, R., Tozzi, F., Muglia, P., Breen, G., Craig, I. W., Müller-Myhsok, B., Kennedy, J. L., Strauss, J., Vincent, J. B., Lewis, C. M., Farmer, A. E., and McGuffin, P. (2013). Estimating the heritability of reporting stressful life events captured by common genetic variants. *Psychological Medicine*, 43(9):1965–1971.

Purcell, S. M., Wray, N. R., Stone, J. L., Visscher, P. M., O'Donovan, M. C., Sullivan, P. F., Sklar, P., (Leader), S. M. P., Ruderfer, D. M., McQuillin, A., Morris, D. W., O'Dushlaine, C. T., Corvin, A., Holmans, P. A., O'Donovan, M. C., Macgregor, S., Gurling, H., Blackwood, D. H. R., Craddock, N. J., Gill, M., Hultman, C. M., Kirov, G. K., Lichtenstein, P., Muir, W. J., Owen, M. J., Pato, C. N., Scolnick, E. M., Clair, D. S., (Leader), P. S., Williams, N. M., Georgieva, L., Nikolov, I., Norton, N., Williams, H., Toncheva, D., Milanova, V., Thelander, E. F., Sullivan, P., O'Dushlaine, C. T., Kenny, E., Quinn, E. M., Choudhury, K., Datta, S., Pimm, J., Thirumalai, S., Puri, V., Krasucki, R., Lawrence, J., Quested, D., Bass, N., Crombie, C., Fraser, G., Kuan, S. L., Walker, N., McGhee, K. A., Pickard, B., Malloy, P., Maclean, A. W., Beck, M. V., Pato, M. T., Medeiros, H., Middleton, F., Carvalho, C., Morley, C., Fanous, A., Conti, D., Knowles, J. A., Ferreira, C. P., Macedo, A., Azevedo, M. H., Kirby, A. N., Ferreira, M. A. R., Daly, M. J., Chambert, K., Kuruvilla, F., Gabriel, S. B., Ardlie, K., and Moran, J. L. (2009). Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*, 460(7256):748–752.

Räsänen, P., Hakko, H., Isohanni, M., Hodgins, S., Järvelin, M.-R., and Tiihonen, J. (1999). Maternal Smoking During Pregnancy and Risk of Criminal Behavior Among Adult Male Offspring in the Northern Finland 1966 Birth Cohort. *American Journal of Psychiatry*, 156(6):857–862.

Reginsson, G. W., Ingason, A., Euesden, J., Bjornsdottir, G., Olafsson, S., Sigurdsson, E., Oskarsson, H., Tyrfingsson, T., Runarsdottir, V., Hansdottir, I., Steinberg, S., Stefansson, H., Gudbjartsson, D. F., Thorgeirsson, T. E., and Stefansson, K. (2017). Polygenic risk scores for schizophrenia and bipolar disorder associate with addiction. *Addiction Biology*, pages n/a–n/a.

Reichenberg, A., Gross, R., Weiser, M., Bresnahan, M., Silverman, J., Harlap, S., Rabinowitz, J., Shulman, C., Malaspina, D., Lubin, G., et al. (2006). Advancing paternal age and autism. *Archives of general psychiatry*, 63(9):1026–1032.

Richmond, R. C., Timpson, N. J., Felix, J. F., Palmer, T., Gaillard, R., McMahon, G., Smith, G. D., Jaddoe, V. W., and Lawlor, D. A. (2017). Using Genetic Variation to Explore the Causal Effect of Maternal Pregnancy Adiposity on Future Offspring Adiposity: A Mendelian Randomisation Study. *PLOS Medicine*, 14(1):e1002221.

Rietveld, C. A., Esko, T., Davies, G., Pers, T. H., Turley, P., Benyamin, B., Chabris, C. F., Emilsson, V., Johnson, A. D., Lee, J. J., Leeuw, C. d., Marioni, R. E., Medland, S. E., Miller, M. B., Rostapshova, O., Lee, S. J. v. d., Vinkhuyzen, A. A. E., Amin, N., Conley, D., Derringer, J., Duijn, C. M. v., Fehrmann, R., Franke, L., Glaeser, E. L., Hansell, N. K., Hayward, C., Iacono, W. G., Ibrahim-Verbaas, C., Jaddoe, V., Karjalainen, J., Laibson, D., Lichtenstein, P., Liewald, D. C., Magnusson, P. K. E., Martin, N. G., McGue, M., McMahon, G., Pedersen, N. L., Pinker, S., Porteous, D. J., Posthuma, D., Rivadeneira, F., Smith, B. H., Starr, J. M., Tiemeier,

H., Timpson, N. J., Trzaskowski, M., Uitterlinden, A. G., Verhulst, F. C., Ward, M. E., Wright, M. J., Smith, G. D., Deary, I. J., Johannesson, M., Plomin, R., Visscher, P. M., Benjamin, D. J., Cesarini, D., and Koellinger, P. D. (2014). Common genetic variants associated with cognitive performance identified using the proxy-phenotype method. *Proceedings of the National Academy of Sciences*, 111(38):13790–13794.

Ripke, S., O'Dushlaine, C., Chambert, K., Moran, J. L., Kähler, A. K., Akterin, S., Bergen, S. E., Collins, A. L., Crowley, J. J., Fromer, M., Kim, Y., Lee, S. H., Magnusson, P. K. E., Sanchez, N., Stahl, E. A., Williams, S., Wray, N. R., Xia, K., Bettella, F., Borglum, A. D., Bulik-Sullivan, B. K., Cormican, P., Craddock, N., de Leeuw, C., Durmishi, N., Gill, M., Golimbet, V., Hamshere, M. L., Holmans, P., Hougaard, D. M., Kendler, K. S., Lin, K., Morris, D. W., Mors, O., Mortensen, P. B., Neale, B. M., O'Neill, F. A., Owen, M. J., Milovancevic, M. P., Posthuma, D., Powell, J., Richards, A. L., Riley, B. P., Ruderfer, D., Rujescu, D., Sigurdsson, E., Silagadze, T., Smit, A. B., Stefansson, H., Steinberg, S., Suvisaari, J., Tosato, S., Verhage, M., Walters, J. T., Multicenter Genetic Studies of Schizophrenia Consortium, Levinson, D. F., Gejman, P. V., Kendler, K. S., Laurent, C., Mowry, B. J., O'Donovan, M. C., Owen, M. J., Pulver, A. E., Riley, B. P., Schwab, S. G., Wildenauer, D. B., Dudbridge, F., Holmans, P., Shi, J., Albus, M., Alexander, M., Campion, D., Cohen, D., Dikeos, D., Duan, J., Eichhammer, P., Godard, S., Hansen, M., Lerer, F. B., Liang, K.-Y., Maier, W., Mallet, J., Nertney, D. A., Nestadt, G., Norton, N., O'Neill, F. A., Papadimitriou, G. N., Ribble, R., Sanders, A. R., Silverman, J. M., Walsh, D., Williams, N. M., Wormley, B., Psychosis Endophenotypes International Consortium, Arranz, M. J., Bakker, S., Bender, S., Bramon, E., Collier, D., Crespo-Facorro, B., Hall, J., Iyegbe, C., Jablensky, A., Kahn, R. S., Kalaydjieva, L., Lawrie, S., Lewis, C. M., Lin, K., Linszen, D. H., Mata, I., McIntosh, A., Murray, R. M., Ophoff, R. A., Powell, J., Rujescu, D., Van Os, J., Walshe, M., Weisbrod, M., Wiersma, D., Wellcome Trust Case Control Consortium 2, Donnelly, P., Barroso, I., Blackwell, J. M., Bramon, E., Brown, M. A., Casas, J. P., Corvin, A. P., Deloukas, P., Duncanson, A., Jankowski, J., Markus, H. S., Mathew, C. G., Palmer, C. N. A., Plomin, R., Rautanen, A., Sawcer, S. J., Trembath, R. C., Viswanathan, A. C., Wood, N. W., Spencer, C. C. A., Band, G., Bellenguez, C., Freeman, C., Hellenthal, G., Giannoulatou, E., Pirinen, M., Pearson, R. D., Strange, A., Su, Z., Vukcevic, D., Donnelly, P., Langford, C., Hunt, S. E., Edkins, S., Gwilliam, R., Blackburn, H., Bumpstead, S. J., Dronov, S., Gillman, M., Gray, E., Hammond, N., Jayakumar, A., McCann, O. T., Liddle, J., Potter, S. C., Ravindrarajah, R., Ricketts, M., Tashakkori-Ghanbaria, A., Waller, M. J., Weston, P., Widaa, S., Whittaker, P., Barroso, I., Deloukas, P., Mathew, C. G., Blackwell, J. M., Brown, M. A., Corvin, A. P., McCarthy, M. I., Spencer, C. C. A., Bramon, E., Corvin, A. P., O'Donovan, M. C., Stefansson, K., Scolnick, E., Purcell, S., McCarroll, S. A., Sklar, P., Hultman, C. M., and Sullivan, P. F. (2013). Genome-wide association analysis identifies 13 new risk loci for schizophrenia. *Nature Genetics*, 45(10):1150–1159.

Ritland, K. (1996). A Marker-Based Method for Inferences About Quantitative Inheritance in Natural Populations. *Evolution*, 50(3):1062–1073.

Robinson, G. K. (1991). That BLUP is a Good Thing: The Estimation of Random Effects. *Statistical Science*, 6(1):15–32.

Robinson, M. R., Wray, N. R., and Visscher, P. M. (2014). Explaining additional genetic variation in complex traits. *Trends in Genetics*, 30(4):124–132.

Rose, R. J., Broms, U., Korhonen, T., Dick, D. M., and Kaprio, J. (2009). Genetics of Smoking Behavior. In *Handbook of Behavior Genetics*, pages 411–432. Springer, New York, NY. DOI: 10.1007/978-0-387-76727-7_28.

Sandin, S., Schendel, D., Magnusson, P., Hultman, C., Surén, P., Susser, E., Grønborg, T., Gissler, M., Gunnes, N., Gross, R., Henning, M., Bresnahan, M., Sourander, A., Hornig, M., Carter, K., Francis, R., Parner, E., Leonard, H., Rosanoff, M., Stoltenberg, C., and Reichenberg, A. (2016). Autism risk associated with parental age and with increasing difference in age between the parents. *Molecular Psychiatry*, 21(5):693–700.

Shi, H., Kichaev, G., and Pasaniuc, B. (2016a). Contrasting the Genetic Architecture of 30 Complex Traits from Summary Association Data. *The American Journal of Human Genetics*, 99(1):139–153.

Shi, H., Mancuso, N., Spendlove, S., and Pasaniuc, B. (2016b). Local genetic correlation gives insights into the shared genetic architecture of complex traits. *bioRxiv*, page 092668.

Shim, H., Chasman, D. I., Smith, J. D., Mora, S., Ridker, P. M., Nickerson, D. A., Krauss, R. M., and Stephens, M. (2015). A Multivariate Genome-Wide Association Analysis of 10 LDL Subfractions, and Their Response to Statin Treatment, in 1868 Caucasians. *PLOS ONE*, 10(4):e0120758.

Silberg, J. L., Maes, H., and Eaves, L. J. (2010). Genetic and environmental influences on the transmission of parental depression to children's depression and conduct disturbance: an extended Children of Twins study. *Journal of Child Psychology and Psychiatry*, 51(6):734–744.

Sirin, S. R. (2005). Socioeconomic Status and Academic Achievement: A Meta-Analytic Review of Research. *Review of Educational Research*, 75(3):417–453.

Solovieff, N., Cotsapas, C., Lee, P. H., Purcell, S. M., and Smoller, J. W. (2013). Pleiotropy in complex traits: challenges and strategies. *Nature Reviews Genetics*, 14(7):483–495.

Speed, D., Cai, N., the UCLEB Consortium, Johnson, M. R., Nejentsev, S., and Balding, D. J. (2017). Reevaluation of SNP heritability in complex human traits. *Nature Genetics*, advance online publication.

Speed, D., Hemani, G., Johnson, M. R., and Balding, D. J. (2012). Improved Heritability Estimation from Genome-wide SNPs. *The American Journal of Human Genetics*, 91(6):1011–1021.

Stephens, M. (2013). A unified framework for association analysis with multiple related phenotypes. *PloS One*, 8(7):e65245.

Szulkin, R., Whitington, T., Eklund, M., Aly, M., Eeles, R. A., Easton, D., Kote-Jarai, Z., Amin Al Olama, A., Benlloch, S., Muir, K., Giles, G. G., Southey, M. C., Fitzgerald, L. M., Henderson, B. E., Schumacher, F., Haiman, C. A., Schleutker, J., Wahlfors, T., Tammela, T. L., Nordestgaard, B. G., Key, T. J., Travis, R. C., Neal, D. E., Donovan, J. L., Hamdy, F. C., Pharoah, P., Pashayan, N., Khaw, K.-T., Stanford, J. L., Thibodeau, S. N., McDonnell, S. K., Schaid, D. J., Maier, C., Vogel, W., Luedeke, M., Herkommer, K., Kibel, A. S., Cybulski, C., Lubiński, J., Kluźniak, W., Cannon-Albright, L., Brenner, H., Butterbach, K., Stegmaier, C., Park, J. Y., Sellers, T., Lim, H.-Y., Slavov, C., Kaneva, R., Mitev, V., Batra, J., Clements, J. A., BioResource, T. A. P. C., Spurdle, A., Teixeira, M. R., Paulo, P., Maia, S., Pandha, H., Michael, A., Kierzek, A., Consortium, t. P., Gronberg, H., and Wiklund, F. (2015). Prediction of individual genetic risk to prostate cancer using a polygenic score. *The Prostate*, 75(13):1467–1474.

Taylor, C. A., Manganello, J. A., Lee, S. J., and Rice, J. C. (2010). Mothers' Spanking of 3-Year-Old Children and Subsequent Risk of Children's Aggressive Behavior. *Pediatrics*, 125(5):e1057–e1065.

The International HapMap Consortium (2005). A haplotype map of the human genome. *Nature*, 437(7063):1299–1320.

Thomas, S. C., Coltman, D. W., and Pemberton, J. M. (2002). The use of marker-based relationship information to estimate the heritability of body weight in a natural population: a cautionary tale. *Journal of Evolutionary Biology*, 15(1):92–99.

Thompson, R. (1973). The Estimation of Variance and Covariance Components with an Application when Records are Subject to Culling. *Biometrics*, 29(3):527–550.

Tobacco and Genetics Consortium (2010). Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nature Genetics*, 42(5):441–447.

Trzaskowski, M., Harlaar, N., Arden, R., Krapohl, E., Rimfeld, K., McMillan, A., Dale, P. S., and Plomin, R. (2014). Genetic influence on family socioeconomic status and children's intelligence. *Intelligence*, 42:83–88.

Turley, P., Walters, R. K., Maghzian, O., Okbay, A., Lee, J. J., Fontana, M. A., Nguyen-Viet, T. A., Furlotte, N. A., Team, a. R., Consortium, S. S. G. A., Magnusson, P., Oskarsson, S., Johannesson, M., Visscher, P. M., Laibson, D., Cesarini, D., Neale, B., and Benjamin, D. J. (2017). MTAG: Multi-Trait Analysis of GWAS. *bioRxiv*, page 118810.

van Kleunen, M. and Ritland, K. (2004). Predicting evolution of floral traits associated with mating system in a natural plant population. *Journal of Evolutionary Biology*, 17(6):1389–1399.

van Kleunen, M. and Ritland, K. (2005). Estimating Heritabilities and Genetic Correlations with Marker-Based Methods: An Experimental Test in Mimulus guttatus. *Journal of Heredity*, 96(4):368–375.

van Rheenen, W., Shatunov, A., Dekker, A. M., McLaughlin, R. L., Diekstra, F. P., Pulit, S. L., van der Spek, R. A. A., Võsa, U., de Jong, S., Robinson, M. R., Yang, J., Fogh, I., van Doormaal, P. T. C., Tazelaar, G. H. P., Koppers, M., Blokhuis, A. M., Sproviero, W., Jones, A. R., Kenna, K. P., van Eijk, K. R., Harschnitz, O., Schellevis, R. D., Brands, W. J., Medic, J., Menelaou, A., Vajda, A., Ticozzi, N., Lin, K., Rogelj, B., Vrabec, K., Ravnik-Glavač, M., Koritnik, B., Zidar, J., Leonardis, L., Grošelj, L. D., Millecamps, S., Salachas, F., Meininger, V., de Carvalho, M., Pinto, S., Mora, J. S., Rojas-García, R., Polak, M., Chandran, S., Colville, S., Swingler, R., Morrison, K. E., Shaw, P. J., Hardy, J., Orrell, R. W., Pittman, A., Sidle, K., Fratta, P., Malaspina, A., Topp, S., Petri, S., Abdulla, S., Drepper, C., Sendtner, M., Meyer, T., Ophoff, R. A., Staats, K. A., Wiedau-Pazos, M., Lomen-Hoerth, C., Van Deerlin, V. M., Trojanowski, J. Q., Elman, L., McCluskey, L., Basak, A. N., Tunca, C., Hamzeiy, H., Parman, Y., Meitinger, T., Lichtner, P., Radivojkov-Blagojevic, M., Andres, C. R., Maurel, C., Bensimon, G., Landwehrmeyer, B., Brice, A., Payan, C. A. M., Saker-Delye, S., Dürr, A., Wood, N. W., Tittmann, L., Lieb, W., Franke, A., Rietschel, M., Cichon, S., Nöthen, M. M., Amouyel, P., Tzourio, C., Dartigues, J.-F., Uitterlinden, A. G., Rivadeneira, F., Estrada, K., Hofman, A., Curtis, C., Blauw, H. M., van der Kooi, A. J., de Visser, M., Goris, A., Weber, M., Shaw, C. E., Smith, B. N., Pansarasa, O., Cereda, C., Del Bo, R., Comi, G. P., D'Alfonso, S., Bertolin, C., Sorarù, G., Mazzini, L., Pensato, V., Gellera, C., Tiloca, C., Ratti, A., Calvo, A., Moglia, C., Brunetti, M., Arcuti, S., Capozzo, R., Zecca, C., Lunetta, C., Penco, S., Riva, N., Padovani, A., Filosto, M., Muller, B., Stuit, R. J., Blair, I., Zhang, K., McCann, E. P., Fifita, J. A., Nicholson, G. A., Rowe, D. B., Pamphlett, R., Kiernan, M. C., Grosskreutz, J., Witte, O. W., Ringer, T., Prell, T., Stubendorff, B., Kurth, I., Hübner, C. A., Leigh, P. N., Casale, F., Chio, A., Beghi, E., Pupillo, E., Tortelli, R., Logroscino, G., Powell, J., Ludolph, A. C., Weishaupt, J. H., Robberecht, W., Van Damme, P., Franke, L., Pers, T. H., Brown, R. H., Glass, J. D., Landers, J. E., Hardiman, O., Andersen, P. M., Corcia, P., Vourc'h, P., Silani, V., Wray, N. R., Visscher, P. M., de Bakker, P. I. W., van Es, M. A., Pasterkamp, R. J., Lewis, C. M., Breen, G., Al-Chalabi, A., van den Berg, L. H., and Veldink, J. H. (2016). Genome-wide association analyses identify new risk variants and the genetic architecture of amyotrophic lateral sclerosis. *Nature Genetics*, 48(9):1043–1048.

Vassos, E., Di Forti, M., Coleman, J., Iyegbe, C., Prata, D., Euesden, J., O'Reilly, P., Curtis, C., Kolliakou, A., Patel, H., Newhouse, S., Traylor, M., Ajnakina, O., Mondelli, V., Marques, T. R., Gardner-Sood, P., Aitchison, K. J., Powell, J., Atakan, Z., Greenwood, K. E., Smith, S., Ismail, K., Pariante, C., Gaughran, F., Dazzan, P., Markus, H. S., David, A. S., Lewis, C. M., Murray, R. M., and Breen, G. (2017). An Examination of Polygenic Score Risk Prediction in Individuals With First-Episode Psychosis. *Biological Psychiatry*, 81(6):470–477.

Vassy, J. L., Hivert, M.-F., Porneala, B., Dauriz, M., Florez, J. C., Dupuis, J., Siscovick, D. S., Fornage, M., Rasmussen-Torvik, L. J., Bouchard, C., and Meigs, J. B. (2014). Polygenic type 2 diabetes prediction at the limit of common variant detection. *Diabetes*, 63(6):2172–2182.

Victora, C. G., Horta, B. L., de Mola, C. L., Quevedo, L., Pinheiro, R. T., Gigante, D. P., Gonçalves, H., and Barros, F. C. (2015). Association between breastfeeding and intelligence, educational attainment, and income at 30 years of age: a prospective birth cohort study from Brazil. *The Lancet Global Health*, 3(4):e199–e205.

Vilhjalmsson, B., Yang, J., Finucane, H. K., Gusev, A., Lindstrom, S., Ripke, S., Genovese, G., Loh, P.-R., Bhatia, G., Do, R., Hayeck, T., Won, H.-H., Consortium, S. W. G. o. t. P. G., the Discovery, B., Kathiresan, S., Pato, M., Pato, C., Tamimi, R., Stahl, E., Zaitlen, N., Pasaniuc, B., Schierup, M., Jager, P. D., Patsopoulos, N., McCarroll, S. A., Daly, M., Purcell, S., Chasman, D., Neale, B., Goddard, M., Visscher, P. M., Kraft, P., Patterson, N. J., and Price, A. L. (2015). Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores. *bioRxiv*, page 015859.

Vinkhuyzen, A. a. E., Van Der Sluis, S., De Geus, E. J. C., Boomsma, D. I., and Posthuma, D. (2010). Genetic influences on 'environmental' factors. *Genes, Brain and Behavior*, 9(3):276–287.

Visscher, P. M., Brown, M. A., McCarthy, M. I., and Yang, J. (2012). Five years of GWAS discovery. *The American Journal of Human Genetics*, 90(1):7–24.

Visscher, P. M., Macgregor, S., Benyamin, B., Zhu, G., Gordon, S., Medland, S., Hill, W. G., Hottenga, J.-J., Willemsen, G., Boomsma, D. I., Liu, Y.-Z., Deng, H.-W., Montgomery, G. W., and Martin, N. G. (2007). Genome Partitioning of Genetic Variation for Height from 11,214 Sibling Pairs. *The American Journal of Human Genetics*, 81(5):1104–1110.

Visscher, P. M., Medland, S. E., Ferreira, M. A. R., Morley, K. I., Zhu, G., Cornes, B. K., Montgomery, G. W., and Martin, N. G. (2006). Assumption-Free Estimation of Heritability from Genome-Wide Identity-by-Descent Sharing between Full Siblings. *PLOS Genetics*, 2(3):e41.

Visscher, P. M., Wray, N. R., Zhang, Q., Sklar, P., McCarthy, M. I., Brown, M. A., and Yang, J. (2017). 10 Years of GWAS Discovery: Biology, Function, and Translation. *The American Journal of Human Genetics*, 101(1):5–22.

Visscher, P. M. and Yang, J. (2016). A plethora of pleiotropy across complex traits. *Nature Genetics*, 48(7):707–708.

Visscher, P. M., Yang, J., and Goddard, M. E. (2010). A commentary on 'common SNPs explain a large proportion of the heritability for human height' by Yang et al. (2010). *Twin Research and Human Genetics: The Official Journal of the International Society for Twin Studies*, 13(6):517–524.

White, I. R., Blane, D., Morris, J. N., and Mourouga, P. (1999). Educational attainment, deprivation-affluence and self reported health in Britain: a cross sectional study. *Journal of Epidemiology and Community Health*, 53(9):535–541.

White, K. R. (1982). The relation between socioeconomic status and academic achievement. *Psychological Bulletin*, 91(3):461–481.

Wiste, A., Robinson, E. B., Milaneschi, Y., Meier, S., Ripke, S., Clements, C. C., Fitzmaurice, G. M., Rietschel, M., Penninx, B. W., Smoller, J. W., and Perlis, R. H. (2014). Bipolar polygenic loading and bipolar spectrum features in major depressive disorder. *Bipolar Disorders*, 16(6):608–616.

Wood, A. R., Esko, T., Yang, J., Vedantam, S., Pers, T. H., Gustafsson, S., Chu, A. Y., Estrada, K., Luan, J., Kutalik, Z., Amin, N., Buchkovich, M. L., Croteau-Chonka, D. C., Day, F. R., Duan, Y., Fall, T., Fehrmann, R., Ferreira, T., Jackson, A. U., Karjalainen, J., Lo, K. S., Locke, A. E., Mägi, R., Mihailov, E., Porcu, E., Randall, J. C., Scherag, A., Vinkhuyzen, A. A. E., Westra, H.-J., Winkler, T. W., Workalemahu, T., Zhao, J. H., Absher, D., Albrecht, E., Anderson, D., Baron, J., Beekman, M., Demirkan, A., Ehret, G. B., Feenstra, B., Feitosa, M. F., Fischer, K., Fraser, R. M., Goel, A., Gong, J., Justice, A. E., Kanoni, S., Kleber, M. E., Kristiansson, K., Lim, U., Lotay, V., Lui, J. C., Mangino, M., Leach, I. M., Medina-Gomez, C., Nalls, M. A., Nyholt, D. R., Palmer, C. D., Pasko, D., Pechlivanis, S., Prokopenko, I., Ried, J. S., Ripke, S., Shungin, D., Stancáková, A., Strawbridge, R. J., Sung, Y. J., Tanaka, T., Teumer, A., Trompet, S., van der Laan, S. W., van Setten, J., Van Vliet-Ostaptchouk, J. V., Wang, Z., Yengo, L., Zhang, W., Afzal, U., Ärnlöv, J., Arscott, G. M., Bandinelli, S., Barrett, A., Bellis, C., Bennett, A. J., Berne, C., Blüher, M., Bolton, J. L., Böttcher, Y., Boyd, H. A., Bruinenberg, M., Buckley, B. M., Buyske, S., Caspersen, I. H., Chines, P. S., Clarke, R., Claudi-Boehm, S., Cooper, M., Daw, E. W., De Jong, P. A., Deelen, J., Delgado, G., Denny, J. C., Dhonukshe-Rutten, R., Dimitriou, M., Doney, A. S. F., Dörr, M., Eklund, N., Eury, E., Folkersen, L., Garcia, M. E., Geller, F., Giedraitis, V., Go, A. S., Grallert, H., Grammer, T. B., Gräßler, J., Grönberg, H., de Groot, L. C. P. G. M., Groves, C. J., Haessler, J., Hall, P., Haller, T., Hallmans, G., Hannemann, A., Hartman, C. A., Hassinen, M., Hayward, C., Heard-Costa, N. L., Helmer, Q., Hemani, G., Henders, A. K., Hillege, H. L., Hlatky, M. A., Hoffmann, W., Hoffmann, P., Holmen, O., Houwing-Duistermaat, J. J., Illig, T., Isaacs, A., James, A. L., Jeff, J., Johansen, B., Johansson, Å., Jolley, J., Juliusdottir, T., Junttila, J., Kho, A. N., Kinnunen, L., Klopp, N., Kocher, T., Kratzer, W., Lichtner, P., Lind, L., Lindström, J., Lobbens, S., Lorentzon, M., Lu, Y., Lyssenko, V., Magnusson, P. K. E., Mahajan, A., Maillard, M., McArdle, W. L., McKenzie, C. A., McLachlan, S., McLaren, P. J., Menni, C., Merger, S., Milani, L., Moayyeri, A., Monda, K. L., Morken, M. A., Müller, G., Müller-Nurasyid, M., Musk, A. W., Narisu, N., Nauck, M., Nolte, I. M., Nöthen, M. M., Oozageer, L., Pilz, S., Rayner, N. W., Renstrom, F., Robertson, N. R., Rose, L. M., Roussel, R., Sanna, S., Scharnagl, H., Scholtens, S., Schumacher, F. R., Schunkert, H., Scott, R. A., Sehmi, J., Seufferlein, T., Shi, J., Silventoinen, K., Smit, J. H., Smith, A. V., Smolonska, J., Stanton, A. V., Stirrups, K., Stott, D. J., Stringham, H. M., Sundström, J., Swertz, M. A., Syvänen, A.-C., Tayo, B. O., Thorleifsson, G., Tyrer, J. P., van Dijk, S., van Schoor, N. M., van der Velde, N., van Heemst, D., van Oort, F. V. A., Vermeulen, S. H., Verweij, N., Vonk, J. M., Waite, L. L., Waldenberger, M., Wennauer, R., Wilkens, L. R., Willenborg, C., Wilsgaard, T., Wojczynski, M. K., Wong, A., Wright, A. F., Zhang, Q., Arveiler, D., Bakker, S. J. L., Beilby, J., Bergman, R. N., Bergmann, S., Biffar, R., Blangero, J., Boomsma, D. I., Bornstein, S. R., Bovet, P., Brambilla, P., Brown, M. J., Campbell, H., Caulfield, M. J., Chakravarti, A., Collins, R., Collins, F. S., Crawford, D. C., Cupples, L. A., Danesh, J., de Faire, U., den Ruijter, H. M., Erbel, R., Erdmann, J., Eriksson, J. G., Farrall, M., Ferrannini, E., Ferrières, J., Ford, I., Forouhi, N. G., Forrester, T., Gansevoort, R. T., Gejman, P. V., Gieger, C., Golay, A., Gottesman, O., Gudnason, V., Gyllensten, U., Haas, D. W., Hall, A. S., Harris, T. B., Hattersley, A. T., Heath, A. C., Hengstenberg, C., Hicks, A. A., Hindorff, L. A., Hingorani, A. D., Hofman, A., Hovingh, G. K., Humphries, S. E., Hunt, S. C., Hypp”onen, E., Jacobs, K. B., Jarvelin, M.-R., Jousilahti, P., Jula, A. M., Kaprio, J., Kastelein, J. J. P., Kayser, M., Kee, F., Keinanen-Kiukaanniemi, S. M., Kiemeney, L. A., Kooner, J. S., Kooperberg,

C., Koskinen, S., Kovacs, P., Kraja, A. T., Kumari, M., Kuusisto, J., Lakka, T. A., Langenberg, C., Le Marchand, L., Lehtimäki, T., Lupoli, S., Madden, P. A. F., Männistö, S., Manunta, P., Marette, A., Matise, T. C., McKnight, B., Meitinger, T., Moll, F. L., Montgomery, G. W., Morris, A. D., Morris, A. P., Murray, J. C., Nelis, M., Ohlsson, C., Oldehinkel, A. J., Ong, K. K., Ouwehand, W. H., Pasterkamp, G., Peters, A., Pramstaller, P. P., Price, J. F., Qi, L., Raitakari, O. T., Rankinen, T., Rao, D. C., Rice, T. K., Ritchie, M., Rudan, I., Salomaa, V., Samani, N. J., Saramies, J., Sarzynski, M. A., Schwarz, P. E. H., Sebert, S., Sever, P., Shuldiner, A. R., Sinisalo, J., Steinthorsdottir, V., Stolk, R. P., Tardif, J.-C., Tönjes, A., Tremblay, A., Tremoli, E., Virtamo, J., Vohl, M.-C., The Electronic Medical Records and Genomics (eMERGE) Consortium, The MIGen Consortium, The PAGE Consortium, The LifeLines Cohort Study, Amouyel, P., Asselbergs, F. W., Assimes, T. L., Bochud, M., Boehm, B. O., Boerwinkle, E., Bottinger, E. P., Bouchard, C., Cauchi, S., Chambers, J. C., Chanock, S. J., Cooper, R. S., de Bakker, P. I. W., Dedoussis, G., Ferrucci, L., Franks, P. W., Froguel, P., Groop, L. C., Haiman, C. A., Hamsten, A., Hayes, M. G., Hui, J., Hunter, D. J., Hveem, K., Jukema, J. W., Kaplan, R. C., Kivimaki, M., Kuh, D., Laakso, M., Liu, Y., Martin, N. G., März, W., Melbye, M., Moebus, S., Munroe, P. B., Njølstad, I., Oostra, B. A., Palmer, C. N. A., Pedersen, N. L., Perola, M., Pérusse, L., Peters, U., Powell, J. E., Power, C., Quertermous, T., Rauramaa, R., Reinmaa, E., Ridker, P. M., Rivadeneira, F., Rotter, J. I., Saaristo, T. E., Saleheen, D., Schlessinger, D., Slagboom, P. E., Snieder, H., Spector, T. D., Strauch, K., Stumvoll, M., Tuomilehto, J., Uusitupa, M., van der Harst, P., Völzke, H., Walker, M., Wareham, N. J., Watkins, H., Wichmann, H.-E., Wilson, J. F., Zanen, P., Deloukas, P., Heid, I. M., Lindgren, C. M., Mohlke, K. L., Speliotes, E. K., Thorsteinsdottir, U., Barroso, I., Fox, C. S., North, K. E., Strachan, D. P., Beckmann, J. S., Berndt, S. I., Boehnke, M., Borecki, I. B., McCarthy, M. I., Metspalu, A., Stefansson, K., Uitterlinden, A. G., van Duijn, C. M., Franke, L., Willer, C. J., Price, A. L., Lettre, G., Loos, R. J. F., Weedon, M. N., Ingelsson, E., O'Connell, J. R., Abecasis, G. R., Chasman, D. I., Goddard, M. E., Visscher, P. M., Hirschhorn, J. N., and Frayling, T. M. (2014). Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature Genetics*, 46(11):1173–1186.

Wray, N. R. (2005). Allele frequencies and the r2 measure of linkage disequilibrium: impact on design and interpretation of association studies. *Twin Research and Human Genetics: The Official Journal of the International Society for Twin Studies*, 8(2):87–94.

Wray, N. R., Goddard, M. E., and Visscher, P. M. (2007). Prediction of individual genetic risk to disease from genome-wide association studies. *Genome Research*, 17(10):1520–1528.

Wray, N. R., Lee, S. H., Mehta, D., Vinkhuyzen, A. A., Dudbridge, F., and Middeldorp, C. M. (2014). Research Review: Polygenic methods and their application to psychiatric traits. *Journal of Child Psychology and Psychiatry*, 55(10):1068–1087.

Wray, N. R., Yang, J., Hayes, B. J., Price, A. L., Goddard, M. E., and Visscher, P. M. (2013). Pitfalls of predicting complex traits from SNPs. *Nature Reviews Genetics*, 14(7):507–515.

Wright, S. (1920). The Relative Importance of Heredity and Environment in Determining the Piebald Pattern of Guinea-Pigs. *Proceedings of the National Academy of Sciences of the United States of America*, 6(6):320–332.

Wright, S. (1984). *Evolution and the Genetics of Populations, Volume 1: Genetic and Biometric Foundations*. University of Chicago Press. Google-Books-ID: 4pTdTWi83ecC.

Yang, J., Bakshi, A., Zhu, Z., Hemani, G., Vinkhuyzen, A. A. E., Lee, S. H., Robinson, M. R., Perry, J. R. B., Nolte, I. M., van Vliet-Ostaptchouk, J. V., Snieder, H., The LifeLines Cohort Study, Esko, T., Milani, L., Mägi, R., Metspalu, A., Hamsten, A., Magnusson, P. K. E., Pedersen,

N. L., Ingelsson, E., Soranzo, N., Keller, M. C., Wray, N. R., Goddard, M. E., and Visscher, P. M. (2015). Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nature Genetics*, advance online publication.

Yang, J., Benyamin, B., McEvoy, B. P., Gordon, S., Henders, A. K., Nyholt, D. R., Madden, P. A., Heath, A. C., Martin, N. G., Montgomery, G. W., Goddard, M. E., and Visscher, P. M. (2010). Common SNPs explain a large proportion of the heritability for human height. *Nature Genetics*, 42(7):565–569.

Yang, J., Lee, S. H., Goddard, M. E., and Visscher, P. M. (2011a). GCTA: A Tool for Genome-wide Complex Trait Analysis. *The American Journal of Human Genetics*, 88(1):76–82.

Yang, J., Lee, S. H., Goddard, M. E., and Visscher, P. M. (2013). Genome-wide complex trait analysis (GCTA): methods, data analyses, and interpretations. In *Genome-Wide Association Studies and Genomic Prediction*, pages 215–236. Springer.

Yang, J., Manolio, T. A., Pasquale, L. R., Boerwinkle, E., Caporaso, N., Cunningham, J. M., de Andrade, M., Feenstra, B., Feingold, E., Hayes, M. G., Hill, W. G., Landi, M. T., Alonso, A., Lettre, G., Lin, P., Ling, H., Lowe, W., Mathias, R. A., Melbye, M., Pugh, E., Cornelis, M. C., Weir, B. S., Goddard, M. E., and Visscher, P. M. (2011b). Genome partitioning of genetic variation for complex traits using common SNPs. *Nature Genetics*, 43(6):519–525.

Zaitlen, N., Kraft, P., Patterson, N., Pasaniuc, B., Bhatia, G., Pollack, S., and Price, A. L. (2013). Using Extended Genealogy to Estimate Components of Heritability for 23 Quantitative and Dichotomous Traits. *PLoS Genet*, 9(5):e1003520.

Zaitlen, N., Pasaniuc, B., Sankararaman, S., Bhatia, G., Zhang, J., Gusev, A., Young, T., Tandon, A., Pollack, S., Vilhjálmsson, B. J., Assimes, T. L., Berndt, S. I., Blot, W. J., Chanock, S., Franceschini, N., Goodman, P. G., He, J., Hennis, A. J. M., Hsing, A., Ingles, S. A., Isaacs, W., Kittles, R. A., Klein, E. A., Lange, L. A., Nemesure, B., Patterson, N., Reich, D., Rybicki, B. A., Stanford, J. L., Stevens, V. L., Strom, S. S., Whitsel, E. A., Witte, J. S., Xu, J., Haiman, C., Wilson, J. G., Kooperberg, C., Stram, D., Reiner, A. P., Tang, H., and Price, A. L. (2014). Leveraging population admixture to characterize the heritability of complex traits. *Nature Genetics*, advance online publication.

Zhang, G., Bacelis, J., Lengyel, C., Teramo, K., Hallman, M., Helgeland, Ø., Johansson, S., Myhre, R., Sengpiel, V., Njølstad, P. R., Jacobsson, B., and Muglia, L. (2015). Assessing the Causal Relationship of Maternal Height on Birth Size and Gestational Age at Birth: A Mendelian Randomization Analysis. *PLOS Medicine*, 12(8):e1001865.

Zheng, J., Erzurumluoglu, A. M., Elsworth, B. L., Kemp, J. P., Howe, L., Haycock, P. C., Hemani, G., Tansey, K., Laurin, C., Pourcain, B. S., Warrington, N. M., Finucane, H. K., Price, A. L., Bulik-Sullivan, B. K., Anttila, V., Paternoster, L., Gaunt, T. R., Evans, D. M., and 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.

Zhou, X. and Stephens, M. (2014). Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nature Methods*, 11(4):407–409.

Zhu, Z., Bakshi, A., Vinkhuyzen, A. A. E., Hemani, G., Lee, S. H., Nolte, I. M., van Vliet-Ostaptchouk, J. V., Snieder, H., Esko, T., Milani, L., Mägi, R., Metspalu, A., Hill, W. G., Weir, B. S., Goddard, M. E., Visscher, P. M., and Yang, J. (2015). Dominance Genetic Variation Contributes Little to the Missing Heritability for Human Complex Traits. *The American Journal of Human Genetics*, 96(3):377–385.