

Gene-environment interactions with essential heterogeneity

GxE with essential heterogeneity

Johannes Hollenbach^{1,2}, Hendrik Schmitz^{1,2}, and Matthias Westphal^{3,1,*}

Abstract: We study how gene-environment interactions between education and genetic endowments affect cognition in old age and use this setting to show that – even with a valid instrument – two-stage least squares (2SLS) estimates of interaction effects can be far away from the true effect. This is the case when treatment effects are heterogeneous and compliance to the instrument depends on the interaction variable. We suggest estimating marginal treatment effects to address this problem. Our estimation results show complementarities between education and genetic predisposition in determining later-life memory. The marginal treatment effect estimates suggest substantially larger gene-environment interactions than the 2SLS estimates.

Keywords: Two-stage least squares, marginal treatment effects, gene-environment interactions, cognitive decline

***Correspondence address:** matthias.westphal@fernuni-hagen.de

We thank the editor and four anonymous referees, as well as Silvia Barcellos, Pietro Biroli, Leandro Carvalho, Jason Fletcher, Andreas Fischeneder, Anna Krumme, Lauren Schmitz, Wendelin Schnedler, and Kevin Thom for their excellent comments and suggestions, which substantially improved the paper. We are also grateful to Souvik Banerjee, Martin Fischer, Hendrik Jürges, and Kristina Strohmaier for their thoughtful feedback as discussants. We further thank participants of the Initiative in Social Genomics group meetings at UW–Madison, the Applied Micro Workshop at UW–Milwaukee, the CINCH-dggö Academy in Health Economics in Essen, the EuHEA PhD Conferences in Bologna and Lucerne, as well as the Brown Bag Seminars in Wuppertal and at RWI Essen for their comments, questions, and lively exchange.

1 Introduction

The recent availability of genetic data has revived the old debate in the social sciences about nature versus nurture in determining success over the life course (see, e.g., [Behrman and Taubman, 1989](#); [Plug and Vijverberg, 2003](#); [Björklund and Salvanes, 2011](#)).

The current focus is on complementarities between both, that is, on estimating gene-environment ($G \times E$) interactions to assess how the effects of environmental exposures or individual decisions vary by genetic endowment. These interaction models are typically specified as

$$Y_i = \beta_0 + \beta_1 E_i + \beta_2 G_i + \beta_3 G_i \times E_i + X_i' \gamma + \varepsilon_i, \quad (1)$$

where Y_i denotes a (long-run) outcome of interest, E_i represents an endogenous environmental exposure or individual decision, G_i captures a pre-determined genetic endowment, and X_i is a vector of control variables. Recent studies have focused on the causal identification of β_1 , β_2 , and β_3 by instrumenting E_i and $G_i \times E_i$ and by removing factors correlated with the environment from G_i (see, e.g., [Schmitz and Conley, 2017](#); [Barcellos *et al.*, 2018, 2021](#); [Pereira *et al.*, 2025](#); [Barcellos *et al.*, 2025](#)). As an alternative to instrumenting E_i , some studies directly estimate interactions of G_i with a plausibly exogenous variable, which we refer to as reduced-form interactions ([van den Berg *et al.*, 2023a,b](#); [Ahlskog *et al.*, 2024](#)).

The focus of our paper is on the estimation of the interaction coefficient, β_3 , which is the central parameter in the gene-environment literature. In its intended interpretation, it measures how the causal effect of the environment varies with genetic endowment, all else being equal. However, as we demonstrate, the commonly used two-stage least squares (2SLS) or reduced-form approaches may not provide reliable estimates of this effect, even with a

valid instrumental variable. This is the case when two conditions hold simultaneously: First, compliers to the instrument for E_i —that is, those who take the treatment E_i only because they are exposed to an instrument Z_i —have different unobserved characteristics between different values of G_i . As an example, among those with a low genetic endowment for education, the share of compliers to a schooling reform may be very different compared to the group with a low genetic endowment for education. Second, the (individual) treatment effects of E_i on Y_i exhibit essential heterogeneity. This occurs when the propensity to take the treatment correlates with the unobserved effect heterogeneity (Heckman *et al.*, 2006). A prominent example of essential heterogeneity is self-selection into treatment based on unobserved gains. For example, those who will benefit most from education are more likely to choose more education. These two conditions frequently occur in real-world settings that are investigated with causal methods. As a result, 2SLS often conflates two different changes when estimating the $G_i \times E_i$ coefficient: first, how the local average treatment effect (LATE) of E_i on Y_i changes with G_i , which is the interaction effect of interest. Second, how the complier subpopulation of this LATE shifts as G_i varies.

In this paper, we (1) comprehensively describe the problem, (2) propose a solution, and (3) apply it to a real-world setting. Using a numerical example, we show that relying on 2SLS estimates of β_3 to provide evidence on how genes and the environment interact can be misleading in a setting with essential heterogeneity and a substantial gradient in the first-stage coefficients across different values of G_i . In our simulation example, the 2SLS coefficient even has the opposite sign of the actual interaction effect. We propose a solution that maintains a fixed underlying population when comparing the effect of E_i on

Y_i for different values of G_i . Estimating marginal treatment effects (MTEs) offers a suitable approach to achieve this (Heckman and Vytlačil, 2005).

We apply this method to the long-term effect of education E_i on cognition in later life Y_i using data from the English Longitudinal Study of Aging (ELSA). We select our sample around the pivotal cohort of a compulsory schooling reform, which extended the minimum school-leaving age from 14 to 15 for individuals born after 1933. Our measure of cognition is the word recall test, a widely used indicator that has been shown to predict cognitive decline (Apolinario *et al.*, 2016; Tsoi *et al.*, 2017). We use data from six waves between 2002 and 2012 when individuals in our sample were between 65 and 80 years old. To measure genetic endowment G_i , we use a polygenic index (PGI), a summary measure that predicts individual-level educational attainment based on the aggregated effects of many DNA differences between individuals. When estimating MTEs, we rely on a recently developed partial identification method by Mogstad *et al.* (2018), also used by Rose and Shem-Tov (2021).

Our paper makes three main contributions to the literature. The first is purely pedagogical. While it is well-documented that selection into gains poses problems for 2SLS in general (see, e.g. Heckman and Vytlačil, 2005), the problem that interaction effects are difficult to interpret in this setting still deserves attention. We aim to provide an accessible and intuitive presentation of the problem. The problem we describe is not limited to the gene-environment literature and is, in principle, relevant to any interaction effect between an endogenous (and instrumented) treatment variable and observable characteristics. In Appendix A we cover other settings where researchers are interested in effect heterogeneity by observables and where the same problems might occur, asking for the same kind of solution.

A more important—and second—contribution is to provide a transparent and easy-to-implement solution by using marginal treatment effects in this setting. Of course, other ways exist to separate possibly correlated observed and unobserved effect heterogeneity. [Kline and Walters \(2019\)](#) discuss the general equivalence between instrumental variable methods and control functions ([Blundell and Powell, 2003](#); [Imbens and Newey, 2009](#)). In control function approaches, the essential heterogeneity is absorbed by a control variable (which might incorporate instrumental variables and functional form assumptions). Although certainly possible, we are not aware of a control function approach in an interaction setting that also uses exogenous variation from instruments. [Arold *et al.* \(2025\)](#) use a control function approach in a $G_i \times E_i$ -study, but they employ the approach by [Altonji and Mansfield \(2018\)](#) and base their control function on group-level averages of observed characteristics without using instruments. Setting up a control function requires the assumption that it is modelled correctly (which may be a stronger assumption than the restrictions we impose). The inability of 2SLS to incorporate unobservable differences between complier groups that could (partially) explain gene-environment interactions is also mentioned in [Barcellos *et al.* \(2021\)](#). They find differences in returns to schooling between individuals with different genetic endowments and use a linear MTE estimation to check whether unobservable factors can explain these disparities, which is not the case in their study.

Our third contribution is a substantive one to the literature on gene-environment interactions, a dynamic field with numerous recent papers in areas related to ours. We are unaware of any study estimating the causal effects of education and its interaction with genetic makeup on memory in later life. [Banks and Mazzonna \(2012\)](#) study the effect of the same reform we do on memory, but without looking at gene-environment interactions.

[Ding *et al.* \(2019\)](#) study the relationship between genes/educational attainment and word recall using data from the Health and Retirement Study (HRS), but do not use exogenous variation in education. [Anderson *et al.* \(2020\)](#) estimate a positive bidirectional relationship between educational attainment and intelligence using genetic variants as instruments. [Schmitz and Conley \(2017\)](#) study whether the effect of the Vietnam War draft lottery on schooling outcomes differs by a genetic predisposition for education. We claim to provide the first evidence of the effects of gene-education interaction on a measure of cognition.

Going beyond genetic gradients in education, [Barcellos *et al.* \(2018\)](#) estimate whether genetic predisposition to obesity moderates the effect of education on health using the UK compulsory schooling reform for the 1957 birth cohort as an exogenous variation and a different data set. [Ahlskog *et al.* \(2024\)](#) estimate reduced-form interactions between compulsory schooling exposure Z_i in Sweden and a set of different PGIs (i.e., they focus on $G_i \times Z_i$) on different outcomes. They find significant interactions for two outcomes (wages and educational attainment), both with the PGIs for educational attainment that we also use in this paper. However, as outlined in Appendix C, one drawback of the focus on $G_i \times Z_i$ is that reduced-form regressions leave the first-stage E_i implicit and do not differentiate between a first-stage gradient in G_i and heterogeneous direct effects of Z_i on the outcome along G_i . Besides [Schmitz and Conley \(2017\)](#) and [Barcellos *et al.* \(2018\)](#), the earliest study in economics on how education can compensate for the effects of genetic differences is probably [Papageorge and Thom \(2020\)](#), who study the impact on labour market outcomes.

Our results are as follows: Applying a benchmark 2SLS estimator, we find a zero effect of education on word recall for individuals in the lowest quintile of G_i —measured by the PGI for education—that is, for those with the lowest genetic propensity for education. On

average, moving to a higher quintile of G_i goes along with an increase in the effect of E_i on Y_i by an insignificant 0.1 words correctly recalled. Using marginal treatment effects, the interaction effect is much larger than when estimating with 2SLS: The average linearized interaction effect across all quintiles of G_i indicates that the effect of E_i on Y_i increases by 0.46–0.47 words per quintile. This corresponds to roughly 10–15 % of the standard deviation of the outcome variable. While education does not improve memory in the group with the lowest genetic endowment, it increases word recall by about 1.8 words in the highest quintile of G_i compared to the lowest. 2SLS would considerably underestimate this gene-environment complementarity. Overall, the 1947 UK compulsory schooling reform has increased schooling, especially for those with lower genetic propensity for schooling. However, these individuals have no returns to schooling in terms of cognition. Instead, significant returns are seen for those with a higher genetic propensity.

The paper proceeds as follows: Section 2 describes the institutional setting of our application and the data used. Section 3 presents 2SLS estimates of gene-environment interactions in our application. Section 4 outlines the challenges in identifying the gene-environment interplay from an econometric perspective and presents our suggested solution. Section 5 gives an overview of the partial identification approach to estimate MTEs and presents our main results. Section 6 concludes.

2 Institutional Setting and Data

2.1 Compulsory schooling reform in the UK

In our application, we exploit exogenous variation from a compulsory schooling reform in the UK. Based on the Education Act of 1944, two reforms were introduced to raise the minimum school-leaving age in England, Scotland, and Wales. We use the first reform, which took effect on April 1, 1947.¹ This reform raised the minimum age for leaving school from 14 to 15. Given that students in the UK typically entered school at age 5, the 1947 reform effectively extended compulsory education from 9 to 10 years. The first birth cohorts to be affected by this change, i.e., the first to be required to attend school for an additional year (the “pivotal cohort”), were those born in April 1933. This particular reform from 1947 has served as exogenous variation for compulsory schooling in studies on the effect of education on wages (Harmon and Walker, 1995; Oreopoulos, 2006; Devereux and Hart, 2010; Clark and Royer, 2013), other labour market outcomes (Clark, 2023), health (Silles, 2009; Powdthavee, 2010; Clark and Royer, 2013; Jürges *et al.*, 2013), health knowledge (Johnston *et al.*, 2015), mortality (Clark and Royer, 2013; Gathmann *et al.*, 2015), and cognitive abilities (Banks and Mazzonna, 2012).

To demonstrate the strong response to the compulsory schooling reform from 1947, Figure 1 shows aggregated cohort-level data from ELSA. It depicts the share of individuals with different levels of schooling by birth cohort. The pivotal cohorts of both compulsory schooling reforms are marked with vertical lines. The highest line (circle markers) shows how the

¹The second reform was enacted much later, in 1972, raising the school-leaving age to 16. Since we are interested in studying memory in old age, we use only the 1947 reform. Cohorts affected by the second reform in 1972 are, for the most part, still too young at the time of data collection for the English Longitudinal Study of Ageing, our data source.

1947 reform caused a significant increase in the share of students leaving school at age 15 or later from about 40% to almost 100%. The middle line (diamond markers) shows how the second reform in 1972 lead to a still remarkable but comparably smaller increase in the share of leaving school at 16 or later from 75% to about 90%. The lowest line (triangle markers) can be read as a placebo test, showing the general trend in increased years of schooling but no discontinuity at the two reform cut-offs ([Clark and Royer, 2013](#)).

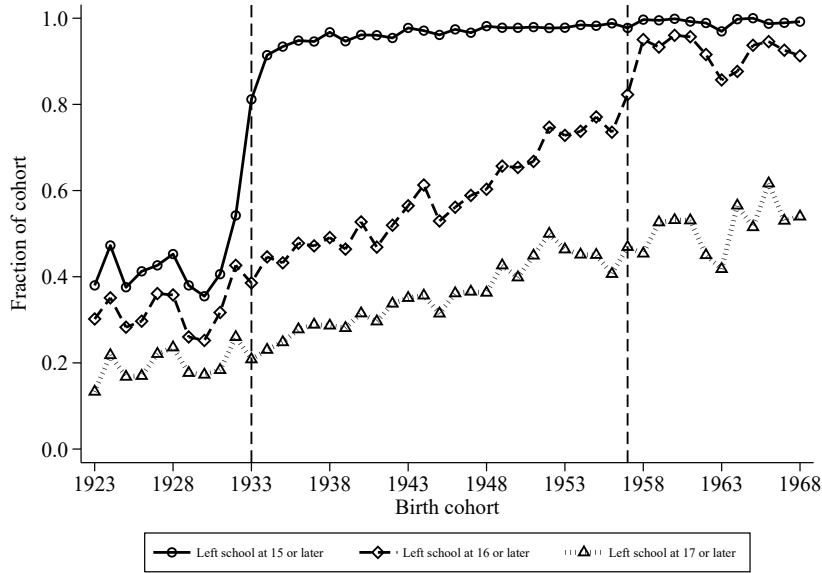


Figure 1: Education by birth cohort

Notes: This figure illustrates the shares of students leaving school at 15 or later, 16 or later, and 17 or later over birth cohorts and how these shares were affected by two compulsory schooling reforms in England using data from ELSA waves 1–9 without the sample selection described in Section 2.2. Vertical dashed lines indicate the first affected birth cohorts of two school-leaving age increases. The three groups are not mutually exclusive and do not add up to 100%. The illustration is adapted from [Clark and Royer \(2013\)](#) to fit our definition of educational attainment.

Besides the high compliance rates, Figure 1 also reveals noncompliance. That is, despite being disallowed by the new compulsory schooling age, some individuals reportedly leave school at the age of 14 after the reform. According to [Clark and Royer \(2013\)](#), who studied

the reforms extensively, this noncompliance is primarily due to individuals born in the summer months who turned 15 after the end of the previous school year but before the start of the next school year.

2.2 Sample and Variables

2.2.1 Sample

We use data from the English Longitudinal Study of Ageing (ELSA), a large representative microdata set providing information on health and other socioeconomic characteristics of individuals aged 50 and over in England ([Banks *et al.*, 2023](#)). ELSA was launched in 2002 and is conducted every two years. It currently comprises eleven waves of interviews.² We use individuals aged 65–80 from waves 1–6 of ELSA. Data collection for wave 6 took place in 2012 and 2013 when individuals born in 1933—our cutoff—turned 80. Thus, starting with wave 7, only individuals born after the cutoff can enter the sample. We exclude the 1933 birth cohort because we lack information on birth month and cannot accurately assign this cohort to pre- or post-reform status (the cutoff is April 1933). We also restrict the data to birth cohorts ten years before and after the reform cut-off.

For our main analysis, we need to limit the data to individuals for whom genetic data is available. This reduces the number of individuals by about 50 % and may introduce a selection bias if the compulsory schooling reform affects the willingness to be genotyped. We find that the sample is selective regarding the outcome variable: Individuals who consent to be genotyped have higher word recall scores on average (see Table [D.1](#) in the Appendix).

²For details of the ELSA sampling procedure, questionnaire content, and fieldwork methodology, see [Steptoe *et al.*, 2013](#).

However, we do not find evidence of a statistically significant effect of the compulsory schooling reform on the probability of being genotyped (see Table D.2 in the Appendix). Similarly, the willingness to be genotyped does not interact with the impact of compulsory schooling on the probability of going to school until at least the age of 15. In our preferred estimation sample, we use all available observations per individual.³ In doing so, we do not assess effects at a single point in time, but implicitly receive average effects over multiple years. The robustness of this choice concerning panel attrition and alternative samples is addressed in Section 3.2. This sample comprises 11,027 observations from 3,009 individuals born between 1923 and 1943, who are observed between 2002 and 2013.

2.2.2 Cognition

Cognitive abilities—as a broad concept—include “the ability to reason, plan, solve problems, think abstractly, comprehend complex ideas, learn quickly and learn from experience. It is not merely book learning, a narrow academic skill, or test-taking smarts. Rather, it reflects a broader and deeper capability for comprehending our surroundings—‘catching on,’ ‘making sense’ of things, or ‘figuring out’ what to do.” (Gottfredson, 1997). The sum of these abilities is called intelligence (Schiele and Schmitz, 2023). A wide range of cognitive tests measure different aspects of cognitive abilities to accommodate this multifaceted notion. ELSA offers several measures, including cognitive capacity, temporal orientation, literacy, and numerical ability. We use test scores from the word recall test in which an interviewer reads ten words to the respondent, who is then asked to recall as many words as possible. This test is administered twice: immediately after the words are read (immediate recall) and

³199 individuals are observed once, 471 twice, 856 three times, 550 four times, 466 five times, and 467 in all six waves.

five minutes later (delayed recall). The scores from both instances are added together to yield a total word recall score, which can range from 0 to 20.

Word recall serves as a measure of episodic memory, susceptible to aging (Rohwedder and Willis, 2010). Episodic memory is considered a component of fluid intelligence, reflecting the innate cognitive ability to store and retrieve information. It is distinct from crystallized intelligence that people acquire over a lifetime (using their fluid intelligence). Word recall has been shown to predict cognitive decline (Bruno *et al.*, 2013; Tsoi *et al.*, 2017) and is an important part of measures for (mild) cognitive impairment (Apolinario *et al.*, 2016; Cadar *et al.*, 2020). Furthermore, it is widely used in economics as a reliable and accessible measure of cognitive functioning (see e.g., Christelis *et al.* 2010; Banks and Mazzonna 2012; Bonsang *et al.* 2012; Schiele and Schmitz 2023). In our estimation sample, the total word recall score, our dependent variable, has a mean of 9.67 correctly recalled words (out of 20) with a standard deviation (SD) of 3.37 words (see Table 1).

2.2.3 Cognition

ELSA does not provide information on respondents' years of education, but on the age at which they completed their continuous full-time education. However, the data is aggregated at the low (finished age 14 or earlier) and high (finished age 19 or later) ends. Our treatment variable E_i is a binary variable equal to one if the individual has left school at 15 or later, and zero otherwise. By design, and as observable in Figure 1, the proportion of individuals having left school at 15 or later (i.e., having stayed in school for at least ten years) is affected by the 1947 education reform that raised the minimum school-leaving age from 14 to 15.

Education is assessed retrospectively, and thus potentially affected by recall bias, a common concern in older age samples. Yet, respondents may better be able to recall the year of school completion (especially so close after the end of World War II) than general years of education. Moreover, Figure 1 (and also the subsequent regression analyses) match remarkably well with the corresponding estimates of [Clark and Royer \(2013\)](#), who use a survey collected from 1991 to 2004—more than one decade before our estimation sample.⁴ Hence, we believe that our education information is unlikely to be significantly affected by recall bias.

2.2.4 Genes

We use an Educational Attainment Polygenic Index (PGI) provided by ELSA and based on [Lee *et al.* \(2018\)](#) to measure genetic makeup. This indicator predicts educational attainment based on differences in genetic variants across individuals. The education PGI we use explains 11-13 % of the variation in educational attainment in the original discovery sample ([Lee *et al.*, 2018](#)). An individual's PGI represents their genetic propensity (or genetic risk—depending on the application) for a particular trait—not just according to one genetic marker, but over many genetic variants. The PGI we are using thus represents individual genetic propensity for educational attainment. For a more detailed explanation of polygenic indices and their construction, see Appendix B. The PGI is normally distributed. Individuals whose genetic endowment puts them on the left side of this distribution have a lower genetic propensity to pursue education; individuals on the right side have a higher propensity. As the choice equations in Section 4 emphasize, this propensity is not deterministic. Individuals with

⁴The corresponding first stage coefficients are 0.445 versus 0.479—a difference smaller than our standard error.

a high PGI are not necessarily highly educated, and highly educated individuals do not necessarily have a high PGI for educational attainment. Our analysis uses the quintiles of this index, yielding five equally sized groups.

When estimating gene-environment interactions, researchers often use a polygenic index of the outcome they are investigating. However, the choice is not set in stone. As [Biroli *et al.* \(2025\)](#) point out, “any PGI could be used if warranted by theory or for empirical reasons”. We target the PGI towards the environmental variable (education) by using an education PGI, and the outcome we investigate is memory. Education PGIs are associated with several different outcomes besides educational attainment: wealth at retirement ([Barth *et al.*, 2020](#)), labour market earnings ([Papageorge and Thom, 2020](#)) and socioeconomic success ([Belsky *et al.*, 2018](#)). In our setting, we can use the education PGI to demonstrate heterogeneous responses to the education reform by the relevant part of the genetic endowment. At the same time, the effect of education on cognition likely varies with genetic propensity for education—which is what we want to estimate.

As control variables in addition to G , we include the first ten principal components of the genetic data, which are summary scores of the overall variation of the genetic data in ELSA, condensed into a smaller number of dimensions. They reflect population stratification, i.e., different frequencies of genetic variants among subpopulations that could be responsible for spurious correlations with outcomes of interest. This could occur if both an outcome and certain genetic variants are more common in one stratum of the population than in another and they do not mate randomly ([Barth *et al.*, 2022](#)). [Price *et al.* \(2006\)](#) show that including principal components as control variables can mitigate these confounding effects. Therefore, adding principal components as controls has become a convention in

gene-environment studies (see, e.g., [Barcellos *et al.*, 2018](#); [Barth *et al.*, 2020, 2022](#); [Pereira *et al.*, 2025](#); [Biroli *et al.*, 2025](#)). Both principal components and polygenic indices combine information from differences of genetic variants across the population. However, they serve different purposes: principal components capture overall genetic similarity and population structure, while PGIs predict specific traits, such as educational attainment, based on gene-outcome associations.⁵

Genes are fixed at conception and, therefore, predetermined. This property is sufficient if one aims to compare causal effects of education across genetic strata (within the same target population). Hence, for our contribution, genes do not have to be exogenous. However, attributing any observed effect heterogeneity solely to genetic factors—as is the ultimate goal in the $G_i \times E_i$ literature—would require identifying genetic impacts within each stratum. Predetermination at conception alone does not guarantee this: unobserved factors may still covary with individuals’ genetic endowments. For example, genetic makeup is inherited across generations, but so too is socio-economic status (for a variety of reasons). Consequently, a correlation between a polygenic index (PGI) and an outcome may reflect the influence of genes, socio-economic background, or a combination of both ([Houmark *et al.*, 2024](#); [Biroli *et al.*, 2025](#)).

We are interested in the effect gradient of E_i (for which we have a well-established instrument) along educational attainment PGI, which aggregates predetermined genetic variants into a single index. Whereas our methodological contribution (identifying the problem and our solution) is unaffected by whether genes are exogenous, we also aim to make a substantive contribution to the $G_i \times E_i$ literature, which discusses the conditions under which

⁵Nevertheless, we show in a robustness check that including principal components does not drive our results (see Table 5).

genes are exogenous. [Houmark *et al.* \(2024\)](#) show that these confounding genetic correlations can be absorbed by observable family characteristics, at least in their data, the Avon Longitudinal Study of Parents and Children. Specifically, controlling for parental education effectively and almost entirely accounts for the correlation not caused by genes. While it remains unclear whether this result also holds for ELSA, we use this insight (and also follow [Barth *et al.*, 2020](#), [Papageorge and Thom, 2020](#), and [Barth *et al.*, 2022](#)) and add parental education measures as additional controls. ELSA includes information on the ages at which a respondent’s mother and father left school, truncated at both ends (ages 14 or under and 19 or over). The vast majority (about 60%) of mothers and fathers left school at age 14 or before. Therefore, we condense the information on education into a categorical variable with three categories: One if both parents have no or low education (i.e., left school at age 14 or earlier), one if at least one parent stayed in school beyond age 14, and one if information on parental education is missing. We have missing information for 988 individuals in our sample. Since we are running local estimations, we choose not to drop them but assign them the category “missing information”. While being a noisy proxy, parental education may capture the family environment also in our setting, allowing us to come closer to isolating the genetic impact in our effect gradient.

2.3 Descriptive Statistics

Table 1 shows descriptive statistics of our main sample of individuals for whom genetic information is available as well as of “treatment” ($E_i = 1$) and “control” ($E_i = 0$) groups separately. Overall, about three-quarters of the observations in the sample are in the treatment group, 66 % were born in 1933 or later, and 52 % are female. The treatment group

scores significantly higher in word recall than the control group. More educated individuals ($E_i = 1$) exhibit a more favourable genetic endowment (significantly less observations in the first and more in the top quintile). Unsurprisingly, individuals in the treatment group are, on average, younger since they are more likely to be born after the compulsory schooling reform. Table D.3 extends the statistics. Table D.4 in the appendix shows the sample means by quintiles of the education PGI. Instrument assignment, age, and proportion of women do not vary across quintiles of the education PGI. However, individuals in higher quintiles perform better on the word recall test. The difference between an average person in the lowest PGI quintile and an average person in the highest quintile is 1.33 words, a sizeable difference compared to the overall mean of 9.67. Not surprisingly, the probability of having more schooling is also higher in higher education PGI quintiles.

3 Benchmark 2SLS estimation

3.1 Empirical Strategy

We start by estimating the gene-environment interactions using “conventional” methods. Since education is a choice variable, an OLS regression will yield biased estimates. We estimate the following 2SLS regression:

$$E_i = \pi_0 + \pi_1 G_i + \pi_2 Z_i + \pi_3 G_i \times Z_i + X' \gamma + f(t) + u_i \quad (2)$$

$$Y_i = \beta_0 + \beta_1 G_i + \beta_2 \widehat{E_i} + \beta_3 \widehat{G_i \times E_i} + X' \delta + f(t) + \varepsilon_i \quad (3)$$

Table 1: Descriptive statistics

	Main sample	By E_i		
	Mean (SD)	$E_i=1$	$E_i=0$	Difference (SE)
<i>Outcome Y_i</i>				
Word recall score	9.67 (3.37)	10.11	8.08	2.03 (0.07)***
<i>Treatment E_i</i>				
Left school ≥ 15	0.78 (0.41)	1.00	0.00	1.00 (0.00)
<i>Polygenic index G_i</i>				
1st PGI quintile	0.20 (0.40)	0.18	0.25	-0.07 (0.01)***
2nd PGI quintile	0.19 (0.40)	0.19	0.21	-0.02 (0.01)**
3rd PGI quintile	0.20 (0.40)	0.21	0.19	0.02 (0.01)**
4th PGI quintile	0.21 (0.41)	0.21	0.20	0.01 (0.01)
5th PGI quintile	0.20 (0.40)	0.22	0.15	0.07 (0.01)***
<i>Instrument Z_i</i>				
Born 1933 or later	0.66 (0.47)	0.82	0.13	0.69 (0.01)***
<i>Selected Controls</i> (for a complete list, see Table D.3)				
Female	0.52 (0.50)	0.52	0.50	0.02 (0.01)**
Birth year	1934.89 (5.00)	1936.29	1929.92	6.37 (0.10)***
Parental education:				
Missing	0.25 (0.43)	0.20	0.41	-0.21 (0.01)***
Both left school ≤ 14	0.57 (0.49)	0.58	0.55	0.03 (0.01)**
At least one left school ≥ 15	0.18 (0.39)	0.22	0.04	0.18 (0.01)***
Observations	11,027	8,590	2,437	

Notes: This table presents descriptive statistics using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. The categories for parental education include: Missing information of at least one parent, both parents left full-time education at age 14 or before or have no education, and at least one parent stayed in school until age 15 or longer. We include mean and standard deviation of the main sample as well as means by E_i , the difference of means and standard errors of a t-test for equality of means. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Eq. (2) is the first stage, where we regress education E_i on our instrument Z_i , genetic predisposition G_i and the interaction of G_i and Z_i .⁶ Eq. (3) shows the second stage. Here, we regress the outcome variable Y_i (the total word recall score for individual i) on the predicted values $\widehat{E_i}$ from the first stage, G_i and the predicted values $\widehat{G_i \times E_i}$. In both stages, we add the same controls X_i which include an indicator variable for sex, the first ten principal components of the genetic data⁷, as well as the ten interactions of the principal components with the instrument, fixed effects for survey wave, as well as a categorical variable for parental education. Furthermore, $f(t)$, is a function that captures a linear cohort trend and its interaction with the instrument Z_i . This specification estimates a fuzzy regression discontinuity model with the re-centred distance to the reform cohort of 1933 (the cohort trend) as the running variable. Finally, u_i and ε_i capture all unobserved factors that affect outcome variables in their respective stages. We cluster standard errors at the individual level in all analyses.

Besides the potential problems due to essential heterogeneity, Eq. (3) linearizes the $G_i \times E_i$ effect (and the effect of G_i itself). This may also mask potentially interesting non-linearities. To be more flexible, we extend our analysis by fully saturating our specification using information on the quintiles of the education polygenic index. These effects compare better to our MTE approach because we directly estimate effects by quintiles. Accordingly, we estimate the following adapted model:

⁶Note that there are technically two first stages, one with the dependent variable E_i and one with the dependent variable $G_i \times E_i$. Depending on how G_i is included, there are more. With G_i as quintiles of the PGI, there are six first stages. For the sake of simplicity, we only show one of them here.

⁷Principal components capture broad genetic patterns in the data, such as similarities among individuals due to shared genetic background. In contrast, genetic predisposition G_i , measured by a PGI, is a prediction of specific traits based on genetic differences between individuals. See Section 2.2 for a detailed description.

$$E_i = \sum_{g=1}^5 \left[\pi_{g,0}^q \mathbb{1}[G_i = g] + \pi_{g,\Delta}^q \mathbb{1}[G_i = g] \times Z_i \right] + X' \gamma^q + f(t) + \omega_i \quad (4)$$

$$Y_i = \sum_{g=1}^5 \beta_{g,0}^q \mathbb{1}[G_i = g] + \beta_{1,1}^q \widehat{E}_i + \sum_{g=2}^5 \beta_{g,1}^q \widehat{\mathbb{1}[G_i = g]} \times E_i + X' \delta^q + f(t) + \eta_i \quad (5)$$

This is the equivalent of the 2SLS model described above in Eqs. (2) and (3), but with sets of indicator variables for the five quintiles of the PGI ($G_i = g$ with $g \in \{1, 2, 3, 4, 5\}$). To distinguish the coefficients from the base model, we add the superscript q . While in the baseline first stage (Eq. 2), π_2 informs about the share of compliers in the data, the $\pi_{1,\Delta}^q$ to $\pi_{5,\Delta}^q$ of Eq. (4) inform about the share of compliers by PGI quintile. In the second stage (Eq. 5), we include \widehat{E}_i as the reference category that captures the local average treatment effect for the lowest quintile ($\beta_{1,1}^q$). The coefficients $\beta_{2,1}^q$ to $\beta_{5,1}^q$ inform about gene-environment interactions relative to the lowest quintile.

3.2 Assumptions

We need to assume that the compulsory schooling reform is a valid instrument for identifying the causal effects of extending schooling beyond age 14. Specifically, the 1933 birth cohort cutoff must be exogenous to the individuals in our sample. This is plausible given that the reform was announced in 1944 and the sample does not suffer from selective non-response or attrition (discussed below). Additionally, we assume that only compulsory schooling changes discontinuously for individuals born after April 1933, without other factors changing simultaneously (the exclusion restriction). Finally, we assume that no individual leaves school earlier because of the reform (the monotonicity assumption).

The exclusion restriction deserves the most discussion, for instance, as spillovers might exist and because two significant events occurred around the time our sample cohorts were born: the Great Depression and World War II. Although individuals may have experienced rationing or evacuations, those on either side of the 1933 cutoff were affected similarly (Clark and Royer, 2013). Furthermore, the compulsory schooling reform may also have increased the general quality of schooling, affecting not only compliers, but also generating spillover effects to always-takers. However, as Clark (2023) documents for the 1947 UK reform, nearly all compliers attended lower-track schools that ended at the minimum leaving age. This makes it unlikely that spillovers to non-complying groups exist. The lower-track schools emphasized practical education, exhibited lower quality (e.g., class size and teacher qualifications), which did not change due to the reform (as resources adjusted to increased enrolment, see Clark and Royer, 2013). These facts suggest that the reform did not affect the quality of schooling (not even for compliers). Thus, we can interpret our treatment effects in terms of years of schooling (as is commonly done in this literature). Clark (2023) also finds that the reform did not raise the probability of students receiving formal academic or vocational qualifications. Nevertheless, as Clark and Royer (2013) note, citing official reports from the period, “the extra year created by the 1947 change introduced some students to more advanced materials and helped other students master more basic material,” suggesting a natural progression in curricula rather than an overhaul.

Additionally, panel attrition may be a concern in older-age samples. If the education reform affected survival or survey participation, it could lead to a disproportionate representation of healthier, more educated individuals among respondents. This selective attrition could bias the estimates if the instrument indirectly influences the sample’s composition

through differential attrition at older ages. [Clark and Royer \(2013\)](#) comprehensively investigated the effect of the 1947 reform on mortality and reported no or negligible effects. Nevertheless, we test for differential panel attrition in our sample by filling in the observations for each individual where necessary from the first wave they were observed in until wave 6, and creating an attrition indicator if they did not respond (for whatever reason) in a subsequent wave. We then regress this indicator on the instrument to assess whether the compulsory schooling reform predicts survey non-response. Table [D.5](#) presents the results. The estimate is negligibly small and not statistically significant. One major difference of our sample compared to related studies (e.g., [Banks and Mazzonna, 2012](#)) is that we exclude individuals who did not provide genetic information to ELSA. Therefore, we also examine whether the probability of sharing genetic information jumps discontinuously at the cutoff (see Table [D.2](#)). We find that this is not the case. We conclude that, although panel attrition is generally a concern, it is not related to the schooling reform in our sample.

3.3 Results

Table [2](#) presents the OLS, reduced form, and 2SLS regression results (Eq. [5](#)) in Columns (1), (2), and (3), respectively. Panel A includes controls for each G_i quintile—but no interaction of G_i with E_i or Z_i . Panel B adds these interactions. Finally, we use the standardized PGI as a continuous interaction variable to show linear effects in Panel C.

3.3.1 OLS

Without interactions, we see a considerable correlation: Individuals who left school at age 15 or later recall about 1.1 words more later in life. Using the compulsory schooling reform

Table 2: OLS, reduced-form and 2SLS estimates

	Dependent variable – total word recall score		
	OLS (1)	Reduced form (2)	2SLS (3)
<i>Panel A: baseline estimate, w/o interactions</i>			
E_i	1.099 (0.138)***		0.154 (0.423)
Z_i		0.075 (0.209)	
<i>Panel B: Including nonlinear interactions with PGI quintiles</i>			
E_i	0.620 (0.245)***		−0.021 (0.449)
Z_i		−0.022 (0.294)	
$G_i = 1$	reference category	reference category	reference category
$G_i = 2$	0.187 (0.269)	0.355 (0.257)	0.239 (0.397)
$G_i = 3$	0.195 (0.274)	0.558 (0.265)**	0.524 (0.451)
$G_i = 4$	0.177 (0.270)	0.802 (0.246)***	0.784 (0.448)*
$G_i = 5$	0.774 (0.314)**	1.084 (0.266)***	0.871 (0.591)
$E_i \times (G_i = 1)$	reference category		reference category
$E_i \times (G_i = 2)$	0.349 (0.325)		0.314 (0.497)
$E_i \times (G_i = 3)$	0.452 (0.331)		0.091 (0.552)
$E_i \times (G_i = 4)$	0.759 (0.328)**		0.049 (0.577)
$E_i \times (G_i = 5)$	0.433 (0.364)		0.394 (0.698)
$Z_i \times (G_i = 1)$		reference category	
$Z_i \times (G_i = 2)$		0.185 (0.317)	
$Z_i \times (G_i = 3)$		0.053 (0.327)	
$Z_i \times (G_i = 4)$		0.028 (0.316)	
$Z_i \times (G_i = 5)$		0.167 (0.330)	
<i>Panel C: Including linear interaction with continuous PGI</i>			
$E_i \times G_i$	0.179 (0.111)		0.011 (0.209)
$Z_i \times G_i$		−0.008 (0.102)	
Controls	Yes	Yes	Yes
Observations	11,027	11,027	11,027

Notes: This table presents OLS, reduced-form and 2SLS estimates of the effect of staying in school until at least age 15 (E_i), an education PGI (G_i), and their gene-environment interaction ($G_i \times E_i$) on word recall later in life using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. In panel A, we show estimates of education (respectively, the instrument Z_i —born in 1933 or later) on word recall without interacting with genetic endowment. For estimates in Panel B, we use quintiles of the polygenic index and estimate non-linear interaction effects. Panel C shows estimates of a linear effect when PGI is standardized and treated as a continuous variable. Coefficients in all panels are obtained from separate regressions. Controls in each case include a linear cohort trend, its interaction with the instrument, sex, survey wave fixed effects, parental education, the first ten principal components of the genetic data as well as interactions of each principal component with the instrument. Standard errors clustered at the individual level shown are in parentheses. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

as exogenous variation, however, the reduced form and the 2SLS estimates suggest that the causal effect, if there is any, is considerably smaller. We find that individuals who would have dropped out at age 14, but had to stay in school for at least one additional year, recall about 0.15 words more. This is a small and statistically insignificant effect. This effect is also more negligible compared to [Banks and Mazzonna \(2012\)](#), who find a relevant and significant impact on the word recall score in their preferred specification. However, as we use an unrestricted sample regarding the school-leaving age, control for G_i (including its principal components) and parental education, and use additional individuals from later waves, our estimates may be more conservative.

The interactions in Panel B use the first quintile ($G_i = 1$) as a reference category, so that the remaining interaction coefficients are interpreted as the additional effect of higher quintiles, relative to the first one. An additional year of education (E_i) is associated with an increase of about 0.62 words later in life for individuals in the lowest PGI quintile. The markups on this association for individuals in the four higher PGI quintiles ($E_i \times (G_i = g)$) are positive across all quintiles, even though their magnitude varies. All in all, this suggests that the association between genes, more education, and memory are mutually reinforcing. When we include the standardized PGI as a continuous variable (in a separate regression, shown in Panel C), its interaction coefficient suggests that a one standard deviation increase in PGI is associated with an additional rise in word recall score by 0.18 words. However, our OLS results only represent correlations.

3.3.2 Reduced-form and IV

Reduced-form estimates that regress word recall on the instrument Z_i and its interaction with G_i are reported in Column (2). Our 2SLS estimates are reported in Column (3). The coefficient of the reduced form without interacting with G_i is almost zero (Panel A). Nevertheless, when considering gene-instrument interactions (Panel B), we see positive effects for all quintiles except the lowest. However, only effects for quintiles two and five have a relevant size, and none of the interaction coefficients are statistically significant. Additionally, we visualize the reduced form alongside the first stage using sample means for each birth cohort in Figure D.1 in the Appendix. The 2SLS regression for education finds a small, positive, but not statistically significant effect of more education on later-life word recall when not interacting with education PGI (Panel A). Furthermore, there is a zero effect of an additional year of schooling on word recall for those in the lowest PGI quintile (Panel B) and a positive estimate for individuals in the upper quintiles. The standard errors are large, so we cannot be certain that these interactions differ from zero. The linear interaction effect using a standardized PGI (Panel C) is also close to zero. Based on these results, we would conclude that, after resolving the endogeneity problem with E_i by instrumenting—if anything—there may only be a small positive interaction effect that cannot be precisely estimated. The cognitive returns to education are likely not much higher for individuals with higher genetic endowment.

However, we outline in Section 4 that the IV estimates are not necessarily consistent estimates of the actual interaction effects.

4 Potential identification problems of interaction effects

4.1 The problem

We are interested in the effect of education E_i on an outcome Y_i and how this effect interacts with genetic endowment G_i . For simplicity, first assume that E_i and G_i are binary variables. Each individual has four potential outcomes, $Y_i^j(G_i = g)$, $j \in \{0, 1\}$, $g \in \{0, 1\}$, where j denotes educational status and G_i the attributes of genetic endowment. For example, $Y_i^1(0)$ is the potential outcome with a high educational level ($E_i = 1$) and a low genetic endowment ($G_i = 0$). However, only one of the four is realized and observed by the researcher. For example, $Y_i^1(1)$ is only observed for an educated person ($E_i = 1$) with good genes ($G_i = 1$). The observation rule is

$$\begin{aligned} Y_i = & E_i \cdot G_i \cdot Y_i^1(1) + E_i \cdot (1 - G_i) \cdot Y_i^1(0) \\ & + (1 - E_i) \cdot G_i \cdot Y_i^0(1) + (1 - E_i) \cdot (1 - G_i) \cdot Y_i^0(0) \end{aligned}$$

Sorting terms, we get

$$\begin{aligned} Y_i = & Y_i^0(0) \\ & + \left(Y_i^1(0) - Y_i^0(0) \right) E_i \\ & + \left(Y_i^0(1) - Y_i^0(0) \right) G_i \\ & + \left(Y_i^1(1) - Y_i^0(1) - [Y_i^1(0) - Y_i^0(0)] \right) G_i \times E_i \end{aligned}$$

This shows how the observation rule corresponds to the $G_i \times E_i$ -workhorse model in Eq. (1):

$$Y_i = \beta_0 + \beta_1 E_i + \beta_2 G_i + \beta_3 G_i \times E_i + \varepsilon_i \quad (6)$$

The gene-environment interaction effect is calculated as $Y_i^1(1) - Y_i^0(1) - (Y_i^1(0) - Y_i^0(0))$, that is, the difference in the effect of E_i on Y_i when $G_i = 1$ (which is $Y_i^1(1) - Y_i^0(1)$) and the effect of E_i on Y_i when $G_i = 0$ (which is $Y_i^1(0) - Y_i^0(0)$).

Assume that G_i is pre-determined while E_i is a choice variable and, therefore, endogenous.⁸ Further assume that we have a binary instrument Z_i that fulfils the classic LATE assumptions (Imbens and Angrist, 1994). For reasons of comparison, recall that, in a model without interactions, the estimated effect of E on Y is

$$\hat{\beta}_1 = \mathbb{E}[Y_i^1 - Y_i^0 | C]$$

where C denotes the compliers, that is, those who take more education if and only if they are affected by the reform. Now, turning back to $G_i \times E_i$ -interaction, we formally show in Appendix E.1 how a 2SLS estimate of β_3 yields

$$\hat{\beta}_3 = \mathbb{E}[Y_i^1(1) - Y_i^0(1) | C(G_i = 1)] - \mathbb{E}[Y_i^1(0) - Y_i^0(0) | C(G_i = 0)] \quad (7)$$

Here, $C(G_i = 1)$ stands for compliers within the group $G_i = 1$ and $C(G_i = 0)$ for compliers within the group $G_i = 0$. This shows that the 2SLS estimate of the interaction coefficient puts together two effects of two potentially different groups: the effect of E_i on Y_i given that $G_i = 1$ in the group $C(G_i = 1)$ and the effect of E_i on Y_i given that $G_i = 0$ in the group

⁸The extension of our framework to an endogenous G_i entails the same kind of problems. Our proposed solution applies to this case but is not straightforward in applications as it requires an instrumental variable for G_i . In Schmitz and Westphal (2025), we apply marginal treatment effect (MTE) estimation with two endogenous variables in a different context, namely causal mediation analysis. However, the estimation of interaction effects with two endogenous variables is beyond the scope of this paper.

$C(G_i = 0)$. Hence, an estimated interaction effect via 2SLS could come from two sources: actual differences in the effect of E_i on Y_i by realization of G_i and/or differences in these effects between the groups $C(G_i = 1)$ and $C(G_i = 0)$.

A simple simulation model visualizes this potential problem. The model and its parametrization are outlined in Appendix E.2. Set up as an illustrative example, Figure 2 shows the average effects of E_i on Y_i (depending on G_i) for four groups in the simulated data. Group 1 on the left are always-takers (AT), irrespective of their realization of G_i . This is because their gains from E_i are so large that they choose more education regardless of Z_i and G_i . The example also produces individuals that are always-takers when $G_i = 1$ but compliers when $G_i = 0$ (Group 2), compliers when $G_i = 1$ and never-takers (NT) when $G_i = 0$ (Group 3) and never-takers, irrespective of G_i (Group 4). Absent simulated data, many of the effects depicted in Figure 2 are unobserved by the researcher.

We sort these four groups along the horizontal axis by their willingness to take education E_i . Those on the left are most willing, and those on the right are least willing. The blue triangles show $\mathbb{E}[Y_i^1(1) - Y_i^0(1)]$, the first part of the interaction effect. The red circles show $\mathbb{E}[Y_i^1(0) - Y_i^0(0)]$, the second part. Thus, the interaction effect for each group is the difference between the blue triangle and the red circle for that group. Our data-generating process is set up so that the interaction effect equals 1.5 for each individual and, consequently, for each group. However, as per Eq. (7), 2SLS calculates it as the difference between the filled blue triangle and the filled red circle, that is $\mathbb{E}[Y_i^1(1) - Y_i^0(1)|C(G_i = 1)] - \mathbb{E}[Y_i^1(0) - Y_i^0(0)|C(G_i = 0)] = 0.2 - 1.4 = -1.2$. Not only is the estimate different in magnitude, but, because of how our example is set up, it is even negative, whereas the true interaction effect is positive.

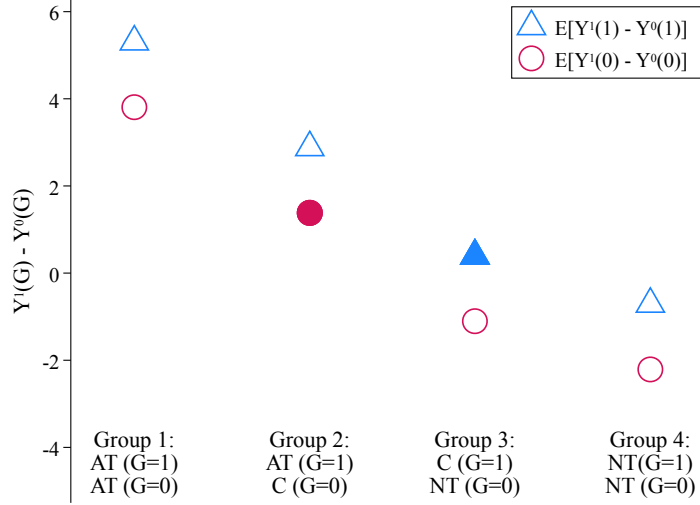


Figure 2: Effects of E_i on Y_i by G_i and complier type in the simulation model

Notes: This figure visualizes stylized potential outcomes from our simulation model. Potential outcomes, their differences and resulting treatment effects are defined by the data-generating process outlined in Appendix E.2 using generated data.

When does this problem occur, and when does it not occur? The figure also helps to grasp the conditions under which 2SLS does *not* fail. First (Condition 1), when the circles and triangles in Figure 2 are on horizontal lines, that is, when both complier groups have the same effects of E_i on Y_i . Technically speaking, this is the case when there is no correlation between willingness to take education—that is, where individuals are located on the x-axis of Figure 2—and the treatment effect—that is, the location on the y-axis. Put differently, when individuals do *not* self-select into the treatment based on their gains from it. A growing literature shows that those who are more likely to take education are also those who benefit more from it (e.g. Carneiro *et al.*, 2011; Nybom, 2017; Kamhöfer *et al.*, 2019). This is often called “selection into gains” or “essential heterogeneity” (Heckman *et al.*, 2006). The simulated data underlying Figure 2 illustrate selection into gains. Those with the highest

effects of education on cognition are those with the highest likelihood to take education. Without this type of selection all red circles in Figure 2 would be on a horizontal line, as would all blue triangles. Then, the effects of E_i on Y_i would not differ by complier type. Even though 2SLS would still make the wrong comparison (filled blue triangle minus filled red circle in Figure 2), the resulting interaction effect would be 1.5.

A second condition (Condition 2) also prevents 2SLS from failing: when the complier groups $C(G_i = 1)$ and $C(G_i = 0)$ do not differ on average. In this case—although there may be differences in potential outcomes between groups—there would be only one complier group, irrespective of G_i and 2SLS would estimate the interaction effect correctly according to Eq. (7). This is the case when G_i does not affect E_i .

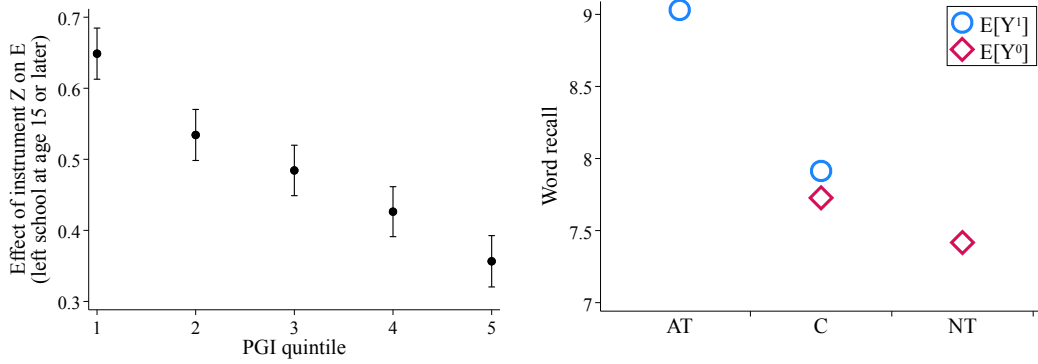
In the simulated data of Figure 2, neither of the two conditions holds. As a result, the 2SLS estimate differs in size from true interaction effects. Turning from the simulated data to our application, we show in the following subsection that neither condition is satisfied in our case.

Before, we note that—as shown in Appendix C—estimating reduced-form $G_i \times Z_i$ interactions does not solve problems with essential heterogeneity because the $G_i \times Z_i$ interaction may be driven solely by a first-stage gradient, even if the (latent structural) $G_i \times E_i$ interaction is absent.

4.2 Evidence suggesting 2SLS may be problematic in our application

We start with Condition 2 and report, in Figure 3a, the results of regressions of E_i on the instrument Z_i , for the five quintiles of G_i . These are the first-stage regressions accompanying

the 2SLS results from Section 3, that is, we report results from Eq. (4). The effects of Z_i on E_i inform about the share of compliers. As reported in Table D.6 in the Appendix, the estimated first-stage coefficients are 0.649 (for $\mathbb{1}[G_i = 1] \times Z_i$), 0.534 (for $\mathbb{1}[G_i = g] \times Z_i$), 0.484 (for $\mathbb{1}[G_i = 3] \times Z_i$), 0.426 (for $\mathbb{1}[G_i = 4] \times Z_i$), and 0.357 (for $\mathbb{1}[G_i = 5] \times Z_i$).



(a) Evidence regarding Condition 2: Effects of Z on E by G (b) Evidence regarding Condition 1: Treated and untreated potential outcomes

Figure 3: Testing for potential problems of 2SLS

Notes: Panel a: This figure shows the complier shares by PGI quintile. These are the coefficients $\pi_{1,\Delta}^q$ to $\pi_{5,\Delta}^q$ from the first-stage Eq. (4). We add 95% confidence intervals. The exact point estimates and their standard errors are also reported in Table D.6. Panel b: Treated and untreated potential outcomes irrespective of G . Red diamonds in refer to potential outcomes for $E_i = 0$; blue circles to $E_i = 1$. The numbers are calculated following the approach by Imbens and Rubin (1997) explained in Appendix E.3. Both panels: data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2.

Overall, we observe a large share of compliers to the education reform in the data. However, it varies substantially across the PGI quintiles. Complier shares monotonically decrease along the PGI. In the lowest quintile ($G_i = 1$), 65 % of all individuals increased their length of education due to the reform. The share of compliers reduces to 36 % in the highest quintile. The compulsory schooling reform had a more substantial impact on individuals in the lowest quintiles of the PGI, who are disadvantaged in terms of the genetic endowment that predicts education. Therefore, the reform was likely effective in targeting

disadvantaged children. It drastically increased their probability of staying in school until at least age 15.

Figure 3b provides some evidence of selection into gains, that is, regarding Condition 1. We use the approach suggested by Imbens and Rubin (1997), to directly estimate $\mathbb{E}[Y_i^1|AT]$, $\mathbb{E}[Y_i^1|C]$, $\mathbb{E}[Y_i^0|C]$, and $\mathbb{E}[Y_i^0|NT]$ from the data, using the LATE assumptions. The formulas behind that are shown in Appendix E.3. For ease of exposition, these numbers do not condition on G . Figure F.1 in the Appendix shows that the findings remain unchanged when we compute these numbers for each quintile of G . Figure 3b shows that the potential outcomes are not horizontal but vary by willingness to take the treatment. Always-takers have larger treated outcomes Y^1 than compliers, while compliers have larger untreated outcomes Y^0 than never-takers. Even though the differences are only small regarding Y^0 , this can be interpreted as evidence of selection into gains.

Selection into gains has been widely documented in the context of education (e.g., Carneiro *et al.*, 2011; Nybom, 2017; Kamhöfer *et al.*, 2019; Westphal *et al.*, 2022; Krumme and Westphal, 2024). Moreover, Barcellos *et al.* (2018) and Barcellos *et al.* (2021) show differences in first-stage responses to a compulsory schooling reform according to G_i . Such self-selection into environments according to genetic makeup has long been established in the $G_i \times E_i$ literature as “active gene-environment correlation”, where the environment mediates the effect of genes on the outcome (Plomin *et al.*, 1977; Plomin, 2014; Biroli *et al.*, 2025).

All in all, our findings provide some evidence that a 2SLS estimation of the $G_i \times E_i$ effect might be problematic in our setting.

4.3 A solution

We suggest going beyond the two points that form the 2SLS estimate, namely the filled circle and triangle in Figure 2. Instead, we propose to estimate the marginal treatment effect (MTE) curve (see, e.g., Heckman and Vytlačil, 2005) by genetic endowment G_i . The MTE framework expands the discrete points from Figure 2 to continuous functions on the unit interval. Here, we provide a verbal account of MTEs and refer to Appendix F.1 for a brief formal introduction, or Heckman and Vytlačil (2005) for an extensive treatment of MTEs.

Figure 4 displays stylized MTE curves from simulated data. On the horizontal axis, instead of grouping individuals as in Figure 2, we now display a continuous index U^E that runs from 0 to 1 and divides the population into percentiles based on their willingness to take education. As in Figure 2, the farther to the left an individual is located—now on the U^E scale—the higher their willingness to take the treatment, and the farther to the right, the lower their willingness. For example, the 10 % of the population with the highest willingness to take education has a $U_i^E \leq 0.1$. Or, the 20% of the population with the highest willingness to take education has a $U_i^E \leq 0.2$. Just as the complier status is unobservable at the individual level, U_i^E is also unobservable at the individual level. In essence, U_i^E measures the most basic building blocks of complying behaviour that are comparable across different covariates (and instruments). As we make comparisons across G_i , it is helpful in our setting, and we can use exogenous variation to identify average characteristics of individuals with specific values of U_i^E .

To define an MTE, compare a local average treatment effect as an aggregated effect for a specific subgroup, e.g. $C(G_i = 1)$, with a marginal treatment effect as an effect defined at

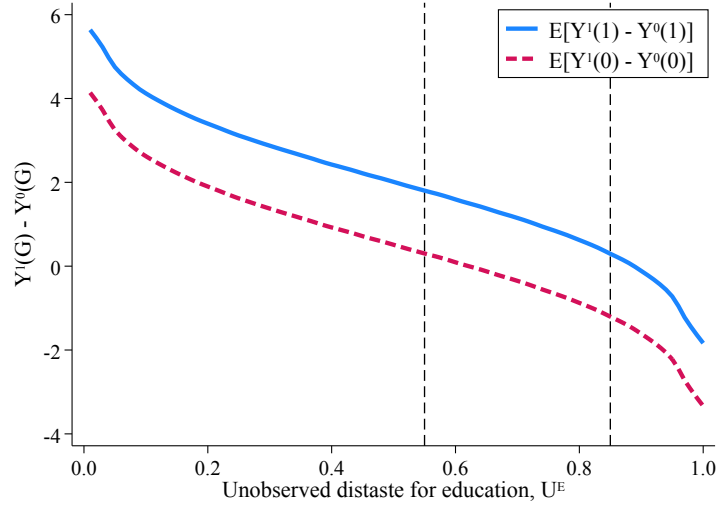


Figure 4: Marginal treatment effects of E_i on Y_i by G_i in the simulation model

Notes: This figure shows stylized marginal treatment effect curves in our simulation model. The differences in potential outcomes are defined by the data-generating process outlined in Appendix E.2 using generated data.

a given value of willingness to take the treatment $U_i^E = u$:

$$LATE(C(G_i = 1), G_i = 1) : \mathbb{E}[Y_i^1(1) - Y_i^0(1) | C(G_i = 1)]$$

$$MTE(u, G_i = 1) : \mathbb{E}[Y_i^1(1) - Y_i^0(1) | U_i^E = u]$$

The MTE is the difference between the two marginal treatment response functions (MTRs): $\mathbb{E}[Y_i^1(1) | U_i^E = u]$ and $\mathbb{E}[Y_i^0(1) | U_i^E = u]$. LATE, MTE, and MTR can likewise be defined for other subgroups, e.g., $G_i = 0$. In general, MTEs can be estimated directly ("joint approach") or by estimating the two MTRs and taking their differences ("separate approach"). In our application below, we will use the separate approach.

Turning back to Figure 4, the interaction effect is the difference between the blue solid and the red dashed curve. In our simulation example, it is always 1.5, but in practice, the two curves do not need to be parallel. The interaction effect may differ along the x-axis, which, as before, represents the willingness to take education.

There are several ways to translate the two curves into interaction effects reflecting the choice of subsamples for whom the effects are estimated: They can be evaluated at certain points on the x-axis or over intervals that represent the location and shares of different groups (always-takers, compliers, never-takers). For example, one possibility is to compute the difference at a specific value of U_i^E , say 0.4. The advantage of this method over 2SLS estimation is that it allows us to make correct vertical comparisons of the two lines and yields a consistent, albeit local, estimate of the interaction at $U_i^E = 0.4$. MTEs can also be used to estimate all treatment parameters, depending on how they are aggregated and how MTEs in different areas of the unit interval are weighted. In principle, it is possible to compute interaction effects using the MTE curves with 2SLS weights either for $C(G_i = 1)$ or $C(G_i = 0)$. In our application below, we use a simpler solution. We will aggregate the MTE results to receive the average interaction effect for all individuals on the U_i^E interval between 0.55 and 0.85, visualized by the two vertical lines in Figure 4. In short, we choose this interval because compliers from every part of the distribution of G_i are located in this area.⁹ This ensures comparability of the MTE result to 2SLS/LATE estimates. We justify this choice in Section 4.5 using Figure 5, which illustrates where compliers are located.

The problem and suggested solution extend beyond gene-environment interactions. They apply to any interaction effect of an endogenous, instrumented treatment with observable

⁹We show robustness checks for other intervals in Table 5.

characteristics, provided there is essential heterogeneity and a first-stage gradient with respect to the interaction variable. Appendix A covers several specific applications from the literature that go beyond the gene-environment setting in which these problems might occur and in which estimating marginal treatment effects could be warranted.

4.4 Going beyond a binary representation of G

The problem and its solution are not specific to cases where G_i is binary. On the one hand, our solution requires a discrete G_i because we will estimate separate curves by G_i . On the other hand, converting a continuous polygenic index into a binary indicator of genetic endowment entails a loss of information. Recall that in our application, we transform the continuous index into quintiles, i.e., a discrete and ordered measure that takes the values $g \in \{1, 2, 3, 4, 5\}$. Consequently, the number of potential outcomes we estimate increases from four to ten. In Table 3, we list these potential outcomes and how to calculate the effect of E_i on Y_i and the $G_i \times E_i$ interaction by genetic type, i.e., quintile of the polygenic index. The reference group is the first (lowest) quintile. Accordingly, all interaction effects are calculated in comparison to this quintile. For example, the gene-environment interaction effect of the fifth (highest) quintile is the difference between the effect of E_i on Y_i for $G = 5$ and the effect of E_i on Y_i for $G = 1$.

Extending the setting to a more complex (but still discrete) classification has advantages. We can make better use of the rich variation of the polygenic index and account for possible nonlinearities in the interaction effects between different sections of the distribution. Of course, the choice to use quintiles is arbitrary. Barcellos *et al.* (2018) and Barcellos *et al.* (2021) show differences in their results according to the terciles of the education polygenic

Table 3: Potential outcomes and calculation of MTEs using quintiles of the polygenic index

$E_i = j$			Individual treatment effects for		
	0	1	the effect of E_i on Y_i	the gene-environment interaction	
$G_i = g$	1	$Y_i^0(1)$	$Y_i^1(1)$	$Y_i^1(1) - Y_i^0(1)$	$(Y_i^1(1) - Y_i^0(1)) - (Y_i^1(1) - Y_i^0(1))$
	2	$Y_i^0(2)$	$Y_i^1(2)$	$Y_i^1(2) - Y_i^0(2)$	$(Y_i^1(2) - Y_i^0(2)) - (Y_i^1(1) - Y_i^0(1))$
	3	$Y_i^0(3)$	$Y_i^1(3)$	$Y_i^1(3) - Y_i^0(3)$	$(Y_i^1(3) - Y_i^0(3)) - (Y_i^1(1) - Y_i^0(1))$
	4	$Y_i^0(4)$	$Y_i^1(4)$	$Y_i^1(4) - Y_i^0(4)$	$(Y_i^1(4) - Y_i^0(4)) - (Y_i^1(1) - Y_i^0(1))$
	5	$Y_i^0(5)$	$Y_i^1(5)$	$Y_i^1(5) - Y_i^0(5)$	$(Y_i^1(5) - Y_i^0(5)) - (Y_i^1(1) - Y_i^0(1))$

Notes: This table lists all combinations of potential outcomes when G_i corresponds to quintiles of the PGI such that $G \in \{1, 2, 3, 4, 5\}$ (left panel). The right panels show how to compute different individual treatment effects, including the interaction effects at every quintile we are after. We chose the first (the lowest) quintile as the reference. All effects are therefore calculated in relation to this group.

index. This is already considerably less restrictive than using a binary representation. The use of quintiles offers a further improvement over terciles. While the general problem is present independent of the PGI's underlying binning, the more granular the binning, the more likely it is to detect the problem by finding a meaningful first-stage gradient or eventual interaction effects. At the same time, it allows us to estimate gene-environment interactions at more points across the polygenic index's distribution, which we can use to detect a possible non-linear evolution of interaction effects across the index. Lastly, using more bins of G_i increases the identifying variation when estimating MTEs with binary instruments.

4.5 Selecting a range for comparable estimates

Individuals respond to the education reform instrument differently by G_i quintile. Figure 3a shows this by highlighting differences in the size of the first stage. We now examine, for each quintile of G_i , the shares and locations on the U^E range of always-takers, compliers, and never-takers. We obtain these using the method of Imbens and Rubin (1997), which we cover in detail in Appendix E.3. Figure 5 visualizes the resulting shares. A large proportion

of those with the lowest genetic propensity for education, i.e., those in $G_i = 1$, are compliers. This is because many of these individuals would have dropped out of school early, but were forced to stay an additional year. Here, compliers are located between U^E values of 0.24 and 0.88. The complier share declines with G_i and is lowest among those with the highest genetic propensity for education ($G_i = 5$), as many of these individuals would have chosen more education regardless of the reform (always-takers). Accordingly, we see the share of always-takers increasing in G_i . Crucially, for $U^E \in [0.55, 0.85]$, we observe compliers, and only compliers, from every G_i quintile. We therefore use this U^E range to obtain MTE estimates that are comparable to, and can be interpreted as, LATEs from a 2SLS estimation.

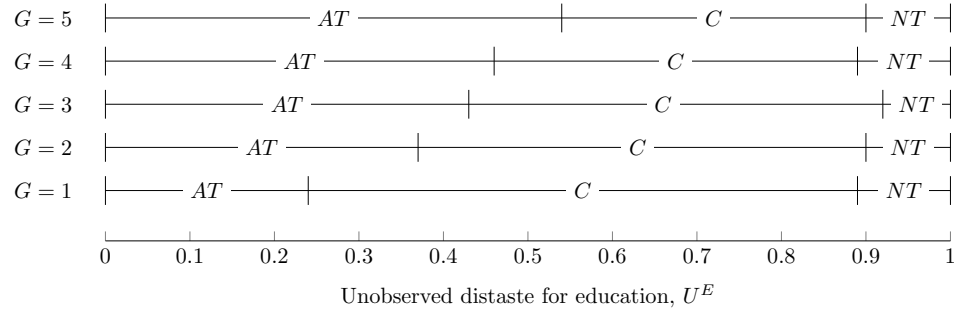


Figure 5: Always-takers, compliers, and never-takers by G

Notes: This figure shows shares and location of always-takers (AT), compliers (C), and never-takers (NT) on the U^E scale in each quintile of G using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. We obtain the shares and thereby their locations using an established method from [Imbens and Rubin \(1997\)](#), which we explain in detail in Section E.3 in the appendix.

5 MTE estimation of the $G_i \times E_i$ interaction

5.1 MTE estimation following Mogstad et al. (2018)

We now describe how, following the partial estimation method suggested by [Mogstad et al. \(2018\)](#), we estimate the ten MTRs needed to calculate the interaction effects. Recall that these ten MTRs are

- $\mathbb{E}[Y_i^1(G_i = 1) \mid U_i^E = u], \quad \mathbb{E}[Y_i^0(G_i = 1) \mid U_i^E = u]$
- ...
- $\mathbb{E}[Y_i^1(G_i = 5) \mid U_i^E = u], \quad \mathbb{E}[Y_i^0(G_i = 5) \mid U_i^E = u]$

These are then used to generate the MTEs as shown in Figure 4.

The method by [Mogstad et al. \(2018\)](#) provides a transparent and credible estimation of MTRs/MTEs when the instrument is binary or when variation in the instrument does not adequately identify MTRs/MTEs across the entire propensity score range. The cost of this method is to give up point identification and resort on bounds of MTEs/MTRs instead. In what follows, we provide a largely verbal overview of the procedure and estimation. Details are available in Appendix F.4. This also includes a step-by-step estimation protocol.

Each of the ten MTRs is approximated by Bernstein polynomials. This means we make parametric functional form assumptions on the MTRs, which, however, are highly flexible and suitable to mimic almost any shape of an MTR on the unit interval one can think of. The question is how to determine the parameters of the Bernstein polynomials. To make this decision, we use a two-step process: First, discard all parameter combinations that lead to MTR shapes that either do not match important observable statistics from

the data or are ruled out by assumptions. Second, among all parameter combinations that fulfill these constraints, we choose those that minimize and those that maximize our target parameter. As a result, we receive bounds on the target parameter. Our target parameter is the difference in the effect of E on Y when we move from $G = 1$ to $G = 5$, evaluated at the U^E -range from 0.55 to 0.85.

In Figure 6, we illustrate the intuition of the method with stylized data. For exposition, we present the procedure as a sequence of steps. In practice, the parameters are found by solving a single optimization problem. Along this example, we present and justify the constraints we make in our application.

CONSTRAINT 1. *All values of $\mathbb{E}[Y_i^j(G_i = g) \mid U_i^E = u, G_i = g]$ lie within the support of Y_i , that is, between 0 and 20.*

This is a trivial constraint. Word recall cannot be below 0 or above 20. Thus, MTRs that produce values of $\mathbb{E}[Y_i^j(G_i = g) \mid U_i^E = u, G_i = g]$ outside this range cannot be suitable. Panel (a) of Figure 6 shows this graphically: Among an enormous number of possible MTRs, we only keep the blue ones that fall between 0 and 20 and we discard the red ones.

CONSTRAINT 2. *When aggregated over the specific interval for the always-takers, compliers, or never-takers, the MTRs reproduce central observable quantities from the data that correspond to these intervals, namely the outcome means $\mathbb{E}[Y_i \mid E_i = 0, Z_i = 0, G_i = g]$ (for NT and C), $\mathbb{E}[Y_i \mid E_i = 0, Z_i = 1, G_i = g]$ (for NT), $\mathbb{E}[Y_i \mid E_i = 1, Z_i = 0, G_i = g]$ (for AT), and $\mathbb{E}[Y_i \mid E_i = 1, Z_i = 1, G_i = g]$ (for AT and C).*

This also implies that they reproduce (now in potential outcome notation) the outcome means for compliers alone $\mathbb{E}[Y_i^1(G_i = g) \mid C(G_i = g)]$, $\mathbb{E}[Y_i^0(G_i = g) \mid C(G_i = g)]$, and the

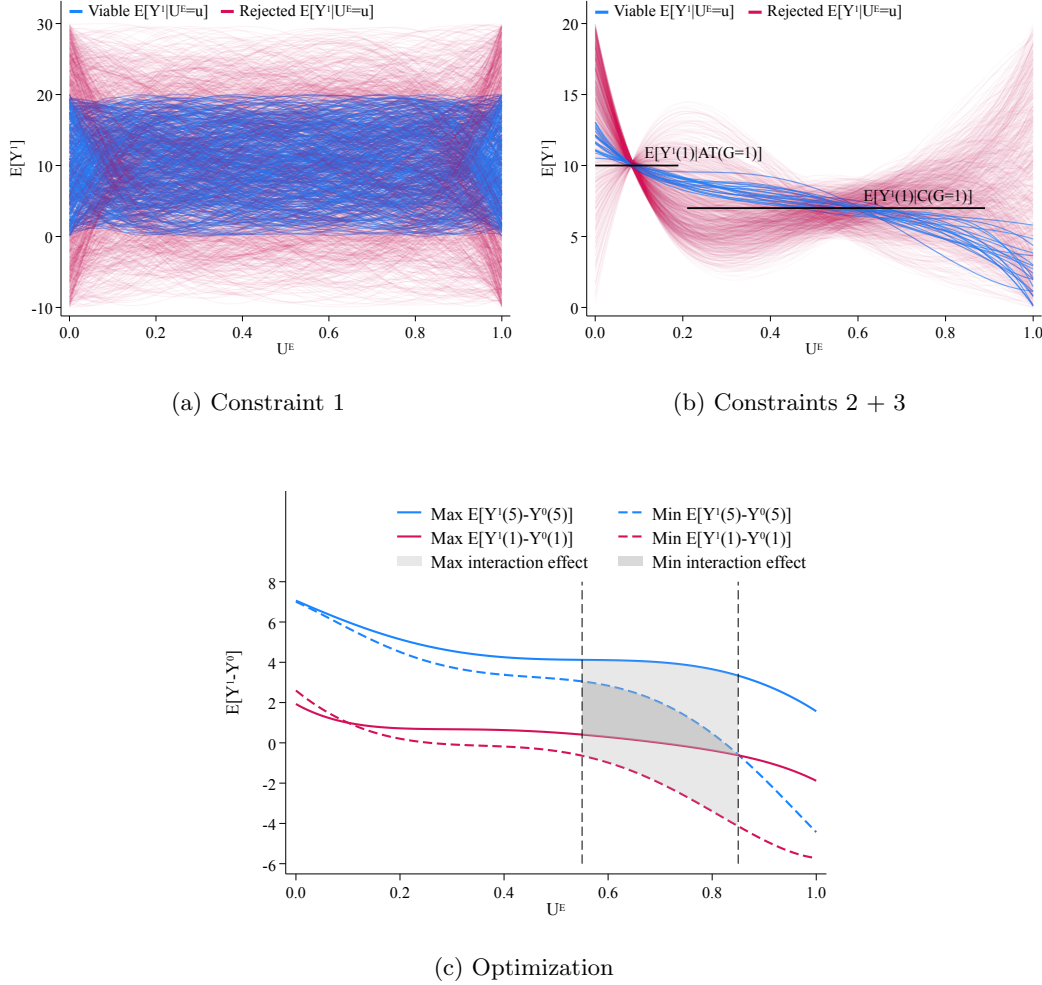


Figure 6: Stylized representation of the constrained optimization of MTE curves

Notes: This figure shows the constrained optimization of the MTE curves to align with the reduced-form evidence and additional restrictions. The graphs are purely illustrative and based on simulated data. Panel (a) shows a large set of MTR curves sharing the same highly flexible parametric structure: Bernstein polynomial curves of degree 5 for Y^1 , without further restrictions. Red lines indicate MTR curves that are incompatible with the support of the dependent variable (word recall score ranging from 0 to 20). In contrast, blue lines represent the remaining curves that lie entirely within this range. Panel (b) introduces two additional constraints. Constraint 2 requires the curve to reproduce the mean recall scores for always-takers, treated compliers, untreated compliers, and never-takers (referred to as the “reduced-form evidence”). All curves whose averages over the U_i^E interval do not match this evidence are discarded. Constraint 3 excludes all MTRs for Y^1 that exhibit a positive slope at any point on the unit interval. Panel (c) displays the final MTE curves that minimize and maximize the interaction effect after applying Constraints 4 and 5. For exposition, we present the estimation as a sequence of steps. In practice, the procedure is solved in a single optimization.

corresponding LATEs for each G_i (Imbens and Rubin, 1997). Note that Constraints 1 and 2 are completely data-driven. We now add three additional constraints derived from assumptions.

CONSTRAINT 3. *Monotone treatment selection* (see Manski, 1997): $\mathbb{E}[Y_i^1(G_i = g)|U_i^E = u]$ does not increase in U_i^E .

Constraint 3 implies—in our application—that among individuals who stay in school beyond age 14 ($E_i = 1$), those with a stronger inclination for education (i.e., with lower U_i^E) do not, on average, have lower word recall ability than those with a weaker preference (higher U_i^E). Our analyses in Figure 3b provides some suggestive evidence in that direction by showing that, in our data, always-takers have higher word recall ability than treated compliers. In Figure F.1 in the appendix, we show that this also holds separately by quintile of G_i . However, we do not impose a comparable constraint on $\mathbb{E}[Y_i^0(G_i = g)|U_i^E = u, G_i = g]$ since the pattern in Figure F.1 is not as clear here.

Panel (b) of Figure 6 combines Constraint 2 and 3. It shows only MTRs that fulfil Constraint 2 and, among them, keeps in blue those that also fulfil Constraint 3.

We proceed by using all viable generated MTRs, that is, $\mathbb{E}[Y_i^0(G)|U_i^E = u, G_i = 1]$ to $\mathbb{E}[Y_i^1(G)|U_i^E = u, G_i = 5]$ to compute the respective MTEs, and retain only those that satisfy Constraints 4 and 5.

CONSTRAINT 4. *No selection into losses*: The MTE $\mathbb{E}[Y_i^1(G_i = g)|U_i^E = u, G_i = g] - \mathbb{E}[Y_i^0(G_i = g)|U_i^E = u, G_i = g]$ is not allowed to increase in U_i^E .

Constraint 4 implies that treatment is a choice and the outcome is beneficial (or correlated with a beneficial variable). This is likely in our setting, where education serves as the

treatment and word recall as the outcome. In such cases, we may expect selection into gains, meaning that MTEs decrease in U_i^E . The literature on the effects of education on earnings and cognitive skills provides extensive empirical evidence supporting selection into gains (Carneiro *et al.*, 2011; Nybom, 2017; Kamhöfer *et al.*, 2019; Westphal *et al.*, 2022; Krumme and Westphal, 2024). Furthermore, in Appendix G, we present suggestive evidence that selection into losses is unlikely in our setting.¹⁰ The constraint also permits MTEs that exhibit no essential heterogeneity, meaning that the MTEs are horizontal. In this case, a 2SLS estimation of $G_i \times E_i$ is unproblematic.

CONSTRAINT 5. *Additive Separability of G_i and the error term: The slope of $\mathbb{E}[Y_i^1(G_i = g)|U_i^E = u, G_i = g]$, $\mathbb{E}[Y_i^0(G_i = g)|U_i^E = u, G_i = g]$, and $\mathbb{E}[Y_i^1(G_i = g)|U_i^E = u, G_i = g] - \mathbb{E}[Y_i^0(G_i = g)|U_i^E = u, G_i = g]$ does not depend on G_i .*

Constraint 5 implies that MTRs and MTEs are parallel across different values of G . While this may be a strong assumption, the ordered and parallel treated potential outcome curves ($E_i = 1$) in Figure F.1 suggest that it can be reasonable. For the untreated outcomes ($E_i = 0$) the pattern is less clear. However, the points may still be consistent with curves that share the same slope across G_i values if nonlinearities along U_i^E are allowed; for example, if MTRs increase up to $U_i^E = 0.9$ and decrease thereafter. Furthermore, by specifying linear regression models such as the workhorse model in Eq. (1), researchers typically impose this assumption implicitly, including in standard 2SLS estimations.

¹⁰This does not mean that all MTE applications find selection into gains. In the context of childcare, Cornelissen *et al.* (2018) find evidence of selection into losses.

Panel (c) of Figure 6 visualizes the final step: among all possible MTRs, we use the ones that maximize (or minimize) our target parameter:

$$\beta_{G \times E}(0.55, 0.85) := \frac{1}{5-1} \int_{0.55}^{0.85} \left[\mathbb{E}[Y_i^1(5) | U_i^E = u] - \mathbb{E}[Y_i^0(5) | U_i^E = u] - \left(\mathbb{E}[Y_i^1(1) | U_i^E = u] - \mathbb{E}[Y_i^0(1) | U_i^E = u] \right) \right] du \quad (8)$$

Equation 8 yields a linearized gene-environment interaction effect on the U_i^E -range that is always covered with compliers from every quintile. The denominator ensures a normalization of the effect to a one-unit increase in G_i . We optimize the interaction effect of the difference between the first and fifth quintile as the natural choice covering the entire distribution of the underlying polygenic index. In Section 5.4, we report robustness checks to show this choice is not crucial.

5.2 Main results

Our main results are visualized in Figure 7. Each panel compares the bounded marginal treatment effects from the first PGI quintile (in red) to the remaining four (in blue). The MTE curves that produce the minimum possible interaction effect are the dashed curves, and the solid curves are MTEs that produce the maximum. Recall that we set up the linear programming approach to optimize the $G_i \times E_i$ effect in the interval $U_i^E \in [0.55, 0.85]$. This is because the compliers from all quintiles are located in this range. In this optimization area, the bounds almost coincide, suggesting that the effects are practically point-identified. This tightness inside the optimization area indicates that the reduced-form evidence (G_i -specific averages for never-takers, always-takers, and compliers) in combination with a highly flexible

polynomial and some additional structure (selection into gains, monotone treatment selection for $E_i = 1$, and additive separability) almost allow for a perfect interpolation of G -specific LATEs to the MTE.¹¹ Note that the tight bounds outside the optimization area are instead a coincidence. The MTEs could look different in this region if the interaction effect were optimized over the whole unit interval. Hence, we only interpret MTE curves and their averages in this region. In Section 5.4, we show that our results are robust to variations of this range. Generally, the differences between the solid MTE curves for quintiles 2–5 and the reference category produce an estimate of the maximal gene-environment interaction effect. The difference between the blue and red dashed curves in each panel yields an estimate of the minimum interaction effect. For example, in the top panel, the area between the blue and red curves indicates how the effect of education on word recall changes in the population when G_i “moves” from the first to the second quintile.

The results have the same sign as our 2SLS estimates. The interaction effect is positive for each quintile comparison. This suggests that individuals with a higher PGI for education benefit more from an additional year of education due to the compulsory schooling reform in terms of their cognition later in life. Our approach also allows us to capture possible nonlinearities in the interaction effect across the PGI. Indeed, the estimated interaction magnitude differs across comparisons. Not surprisingly, the highest quintile has the largest interaction effect. However, the interaction size for the second quintile is substantial. Those in the third and fourth quintiles have the smallest effects.

We present estimates of the nonlinear interaction effects in Panel A of Table 4. While the previously discussed 2SLS results from Table 2 are reported in column (1) as a benchmark,

¹¹Figure D.2 shows the bounds when never-takers are excluded. The shape of the polynomials is different because they are anchored at less points. As a result, the bounds are wider. We discuss this in more detail in Section 5.4.

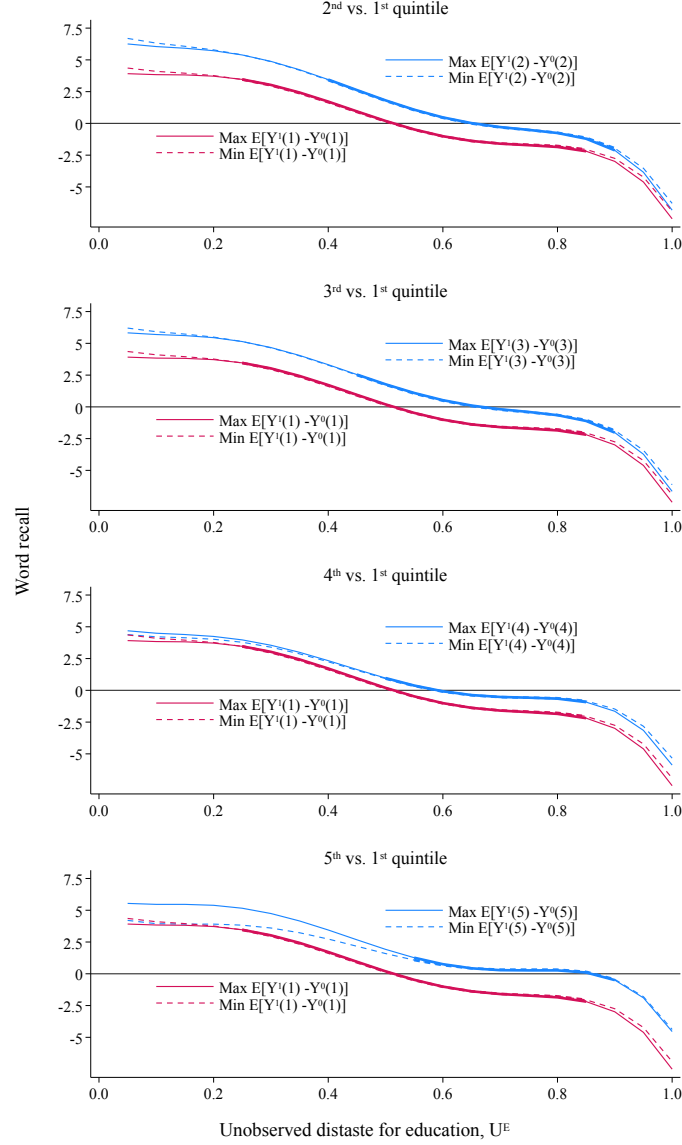


Figure 7: Quintile comparisons of the interaction effect

Notes: This figure shows the four comparisons of gene-environment interactions from our bounding approach using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. For every PGI quintile, we estimate bounds: maxima (solid lines) and minima (dashed lines) at which the interaction effect is maximized/minimized. The bounds for quintiles 2–5 (in blue) are compared to those of the bottom quintile (in red), our reference category, yielding four comparisons. The smallest possible gene-environment interaction is the difference between the blue and red dashed curves over $U_i^E \in [0.55, 0.85]$; the largest possible interaction effect is calculated as the difference between blue and red solid curves over this interval. The thick part of the curves indicates the size of the complier share and its location on the U_i^E scale, both of which differ by PGI quintile.

columns (2) and (3) present the bounds of the marginal treatment effects from Figure 7 aggregated over the U_i^E range from 0.55 to 0.85. As in Table 2, the effect on E_i in the first row indicates the baseline effect in the bottom quintile. The direct effects on G_i in the subsequent rows are not of immediate interest, but we present them for completeness. Our focus is on the interaction effects, which are again relative to the reference category, the bottom quintile.

In addition, we present the linearized interaction effect from the quintile coefficients (Panel B, see Eq. 8) that is our main measure of the gene-environment interaction effect.¹² This measure is simply the slope of a line through the interaction effect estimates of the lowest and highest quintiles and can be thought of as the average interaction effect standardized to a one-quintile change.¹³ This allows for a comparison of interaction effects from 2SLS and MTE in one number to infer whether unobserved effect heterogeneity and different proportions of compliers in G_i —which we fix by estimating marginal treatment effects—affected the 2SLS coefficients.

Overall, four features characterize our results. First, the MTE method yields informative and narrow upper and lower bounds of the interaction MTE, which almost point-identify the effect. Second, even the lower-bound MTE results indicate a relevant interaction effect that is substantially larger than 2SLS estimates. The linearized lower bound is about 4.7 times larger than the linearized 2SLS coefficient. While we could not detect significant gene-environment interaction effects with 2SLS, MTE estimation suggests statistically significant effects at the 5 % level. However, note that the most important difference is the effect sizes

¹²The linear slope is calculated as $(\beta_{5,1}^g - \beta_{1,1}^g)/4$. The interaction coefficient for the bottom quintile, $\beta_{1,1}^g$, is zero since this quintile serves as the reference category.

¹³Note that this measure differs from the linear interaction coefficients in Table 2, Panel C, where we present interactions with the standardized PGI as a continuous variable, which is conceptually different.

Table 4: Estimates of the $G_i \times E_i$ interaction

	Dependent variable – total word recall score					
	2SLS (1)		MTE _{min} (2)		MTE _{max} (3)	
<i>Panel A: nonlinear $G \times E$ effect with G_i as quintiles</i>						
E_i	−0.021	(0.449)	0.121	(0.450)	0.121	(0.450)
$G_i = 1$	reference category		reference category		reference category	
$G_i = 2$	0.239	(0.397)	−0.415	(0.463)	−0.415	(0.441)
$G_i = 3$	0.524	(0.451)	−0.576	(0.607)	−0.579	(0.579)
$G_i = 4$	0.784	(0.448)*	−0.249	(0.479)	−0.252	(0.547)
$G_i = 5$	0.871	(0.591)	0.095	(0.822)	0.078	(0.840)
$E_i \times (G_i = 1)$	reference category		reference category		reference category	
$E_i \times (G_i = 2)$	0.314	(0.497)	1.308	(0.582)**	1.342	(0.766)*
$E_i \times (G_i = 3)$	0.091	(0.552)	1.377	(0.666)**	1.418	(0.852)*
$E_i \times (G_i = 4)$	0.049	(0.577)	1.012	(0.637)	1.033	(0.771)
$E_i \times (G_i = 5)$	0.394	(0.698)	1.851	(0.810)**	1.883	(0.912)**
<i>Panel B: linearized $G_i \times E_i$ effect from quintile coefficients</i>						
$E_i \times G_i$	0.098	(0.174)	0.463	(0.203)**	0.471	(0.228)**
Controls	Yes		Yes		Yes	
Observations	11,027		11,027		11,027	

Notes: This table presents 2SLS and MTE estimates of the effect of staying in school until at least age 15 (E_i), an education PGI (G_i), and their gene-environment interaction ($G_i \times E_i$) on word recall later in life using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. Panel A shows estimates for which we use quintiles of the PGI to estimate possible nonlinear effects across G_i . Estimates that include G_i are computed relative to the reference category, the bottom quintile. Panel B shows the average linearized $G_i \times E_i$ effect across all quintiles of G_i ; our main result. For reference, the average recall score in the sample is 9.67 words with a standard deviation of 3.37. Consequently, the linearized MTE estimates suggest a $G_i \times E_i$ interaction effect of about 14 % of a standard deviation per quintile of G_i . 2SLS estimates from Table 2 are included for reference in Column (1). The MTE estimates in column (2) refer to the minimal effects where the underlying optimization minimizes the linearized interaction effect. Estimates in column (3) are the maximal effects estimated accordingly. The controls in each case include a linear cohort trend, its interaction with the instrument, gender, survey wave fixed effects, parental education, the first ten principal components of the genetic data, and their interactions with the instrument. Results in different panels are obtained from separate regressions. Standard errors clustered at the individual level are shown in parentheses. For MTE bounds, standard errors are bootstrapped with 100 repetitions. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

and not the statistical significance. Because it is based on a two-sided test, the reported significance levels for the bounds are stricter than necessary if the interest lies in the true value within the bound. Imbens and Manski (2004) suggests that a one-sided test is sufficient and would lead all interaction coefficients except for ($E_i \times G_i = 2$) to shift one significance level (i.e., gain one star). The MTE standard errors are only slightly larger than those of the 2SLS. This is because our MTE estimation is only slightly more data demanding than a

pure 2SLS estimation, our target parameter does not rely on extrapolation to non-complying groups, and we additionally impose the outlined restrictions that reduce sampling variation of the MTE relative to 2SLS.¹⁴ Third, our estimates suggest that the gene-environment interaction is more substantial for individuals with higher PGI, while 2SLS estimates suggest a zero or small and statistically insignificant interaction effect. On average, “moving” to a higher PGI quintile leads to an additional increase of 0.46-0.47 words in the impact of compulsory education on word recall due to the education reform. This finding reveals substantial heterogeneity and suggests a high complementarity between education and “nature” as measured by the PGI. Individuals with a higher PGI have higher returns to schooling in terms of cognitive abilities later in life. This result is independent of observable and unobservable factors, both of which we can fix by estimating marginal treatment effects. Fourth, the interaction effect size does not appear linear along the PGI, as indicated by the visual differences in the interaction effects between each panel in Figure 7. The MTE results suggest that individuals in the highest quintile experience a large additional increase in word recall of between 1.85 and 1.88 words relative to those in the first quintile. The increases for individuals in the fourth quintile may not be statistically different from the interaction for individuals in the lowest quintile, although the point estimates are also positive and substantial.

To put our results into perspective, they suggest that individuals in the lowest quintile of the education PGI do not experience increased memory later in life. However, compared to them, additional education increases memory by about half a standard deviation (i.e., 50% of

¹⁴The method of Mogstad and Torgovitsky (2018) only requires estimating twice the number of parameters compared to 2SLS (see Eq. E.9) and without restricting the sample. The method does not extrapolate to non-complying individuals (which could inflate the standard errors), because we report MTE-based estimates for an U_i^E interval that is only covered by compliers in every G quintile.

3.37) for those in the highest quintile. Individuals in between experience lower benefits from schooling than those in the top quintile. The average linearized per-quintile increase across quintiles 2–5 is about 14 % of a standard deviation. We also calculate MTE estimates of the total effect of E_i on Y_i , i.e., without interacting with G_i . The results are reported in Table D.7. The lower-bound MTE estimate is small, positive, and not statistically significant. The upper bound suggest a total effect of 1.64 words, or, again, half of a standard deviation. Conditional on the standard error of the linearized $E_i \times G_i$ effect (0.203) and detecting a significant effect, we have a "minimal significant effect size" of $1.96 \times 0.203 = 0.398$.¹⁵ This effect is 11.8 % of the standard deviation of the unconditional recall error. This shows that our statistical power is too low to detect minor interaction effects. For example, it would be difficult to detect a significant 2SLS interaction given its small point estimate. Nonetheless, Table D.8 documents that our power suffices to detect the average gradient along unobserved heterogeneity as found in the MTE literature.¹⁶

5.3 Discussion

How do our results relate to the previous literature? To begin with, we compare our results with studies that estimate the effects of schooling on cognitive measures without accounting for gene-environment interactions. Related studies find relatively large effects of schooling on word recall ability. Our maximum MTE estimate (without interactions), as well as the estimated interaction effect comparing individuals at the lowest and highest ends of the

¹⁵This minimal significant effect size differs from the minimal detectable effect size for power calculations because we condition on the test result. The minimal detectable effect size is typically defined as the minimal effect detected as significant at the 5 % level in 80 % of all cases. This measure, however, cannot be computed without further assumptions.

¹⁶A comparison between the interquintile heterogeneity along the PGI in our study and the quantiles of the unobserved (MTE) heterogeneity is appropriate because (i) this heterogeneity is supposed to correlate highly (as uncontrolled genes are supposed to be an important component of the unobserved heterogeneity) and (ii) it is easy to transfer the MTE heterogeneity (reported along quantiles) to a measure along quintiles.

education PGI distribution, are within the range of main effects reported in this literature (see the summary of our results above): Using the same setting as in this paper, [Banks and Mazzonna \(2012\)](#) find increases of about half a standard deviation in old-age memory from the additional year of schooling induced by the 1947 UK reform. [Gorman \(2023\)](#) finds increases of one-third to half a standard deviation in memory from the 1972 reform, and [Glymour *et al.* \(2008\)](#) report about a third of a standard deviation for U.S. compulsory schooling increases. [Carvalho \(2025\)](#), among other things, estimates the effect of the 1972 education reform on fluid intelligence, finding no effect. However, his outcome—answers to several reasoning questions—captures a different aspect of fluid intelligence related to problem-solving, whereas our outcome measures episodic memory.

Part of the effect of schooling on improved old-age cognitive abilities could arise through higher earnings. However, while [Harmon and Walker \(1995\)](#) and [Oreopoulos \(2006\)](#) find substantial effects of schooling on wages from the 1947 reform, [Devereux and Hart \(2010\)](#) uncover that the impact on earnings is considerably smaller. Taken together, findings in the literature imply that education could have a stronger influence on cognition than on wages. [Banks and Mazzonna \(2012\)](#) discuss channels that include and extend beyond income. They suggest that the effect of education on cognition could also come about via access to more cognitively demanding occupations, enabling greater engagement in cognitively stimulating activities; potentially increased social and cultural participation, or greater productive efficiency in maintaining cognitive health. They rule out effects via physical health improvements and mortality, and note that benefits likely emerged among lower-educated individuals due to diminishing marginal returns to school years and possibly the protective effect of delayed entry into the labour force.

To our knowledge, we are the first to estimate gene-environment interactions on old-age memory. Nevertheless, it is worthwhile to compare interactions of education and genetic markers on different outcomes found in related studies, especially income. Using the 1972 UK increase in the school-leaving age [Barcellos *et al.* \(2021\)](#) estimate average wage increases of 6-7% and interaction effects of 2% additional gains per standard deviation of an education PGI. When dividing the PGI into terciles, they find no wage effects for individuals with low PGI, and increases of 6-8% for those in the highest PGI tercile. [Ahlskog *et al.* \(2024\)](#) show that a Swedish schooling reform directly benefited earnings of women with lower education PGI, finding interaction effects of 12% of the average reform effect—however, only for women from wealthy families. Our gene-environment interaction estimates for cognition are larger than those for income, which aligns with the general notion that effects of education on cognition are larger than wage effects. In a recent addition, [Barcellos *et al.* \(2025\)](#) estimate a large $G_i \times E_i$ interaction for an Alzheimer’s diagnosis. They show that schooling is especially beneficial for people with higher genetic risk, reducing their likelihood of an Alzheimer’s diagnosis by at least 40% of the pre-reform average.

5.4 Robustness

We perform several robustness checks and report our main measure, the linearized MTE bound estimates, in Table 5. As a baseline for reference, we provide our main result from Table 4. First, we estimate the interaction over alternative U_i^E ranges, in particular $U_i^E \in [0.6, 0.8]$ and $U_i^E \in [0.5, 0.9]$. While the main range $U_i^E \in [0.55, 0.85]$ covers most compliers across all quintiles, we show that this choice is not critical to our main results (see Panel B). Both over a wider and a narrower U_i^E range, the distance between minimum and maximum

bounds only marginally varies and remains statistically significant at the 5% level on a two-sided test. We do not go beyond $U_i^E \in [0.5, 0.9]$ since the never-takers are predominantly located to the right of $U_i^E = 0.9$.

Second, we show robustness checks for different sample compositions. Our dataset consists of repeated cross-sections (waves) of ELSA, and we control for wave fixed effects. Nevertheless, some individuals are observed only once, while most are observed several times. We include a robustness check that uses only the most recent observation for each individual. This reduces the number of observations (from 11,027 to 3,009) but not the number of individuals in the analysis. Both the upper and lower bound estimates are slightly larger, as are the standard errors due to the lower number of observations. Both estimates are statistically significant at the 10% level on the default two-sided test.

Furthermore, we include individuals under 65, indicating that the sample composition, particularly our focus on older individuals, does not significantly alter our results. The minimum MTE estimate is smaller than our main result, but the maximum remains similar to our main specification. The choice of the appropriate polynomial to control for cohort trends is not obvious. We demonstrate how the results change when quadratic cohort trends are used instead of linear ones, or when cohort trends are allowed to vary linearly across PGI quintiles. Doing so increases flexibility. These changes only marginally affect the possible range of effects, as the estimates are similar to those in our main result. However, standard errors increase, especially when interacting cohort trends with G_i . With the latter, the minimum and maximum estimates both increase slightly. With squared trends, however, the minimum decreases and the maximum increases somewhat.

Next, we show that controlling for principal components of the genetic data—while being sensible and an established norm in the literature—does not drive our main result. Removing them and their interactions with the instrument as control variables leads to slightly lower estimates for both bounds. Next, we remove observed never-takers from the analysis (see Section F.3 for a discussion of never-takers). Their presence helps to tighten the bounds. However, even without them, the minimum and maximum MTEs are informative. The interaction effect’s lower and upper bounds are still positive. However, the lower bound may not be statistically different from zero, and the upper bound is larger than in our main result. This is to be expected, since never-takers have lower expected outcomes. Not including them in the analysis means that the MTE bounds do not have to reproduce these lower means. As a result, the resulting MTE curves will look different. We visualize the quintile comparisons when computing interaction effects without never-takers in Figure D.2 in the Appendix.

Next, we consider the delayed word recall score as an alternative outcome variable. This measure counts the number of correctly recalled words (out of ten) five minutes after they are read to participants in studies like ELSA. Since delayed word recall is a more difficult task than recalling the words immediately—the second component of our main outcome, total word recall—the sample mean for delayed word recall is 4.15, less than half the total word recall mean. As expected, the estimated MTE effect sizes are much smaller, even though the lower bound is not statistically significant. Furthermore, we estimate our main result, but calculate standard errors with double the number of bootstrap repetitions (200). The standard errors barely change.

We show additional robustness checks of our linearized 2SLS measure in Table D.9. Given the sample size, we are somewhat limited in how flexibly we can estimate the 2SLS model.

Table 5: Robustness

Linearized $G \times E$ effect	Dependent variable – total word recall score	
	MTE _{min} (1)	MTE _{max} (2)
Baseline (main result, Table 4)	0.463 (0.203)**	0.471 (0.228)**
$U_i^E \in [0.6, 0.8]$	0.461 (0.203)**	0.466 (0.232)**
$U_i^E \in [0.5, 0.9]$	0.418 (0.201)**	0.507 (0.217)**
One observation per individual	0.480 (0.248)*	0.496 (0.277)*
Keeping individuals below age 65	0.355 (0.206)*	0.482 (0.201)**
Squared cohort trends	0.444 (0.210)**	0.512 (0.234)**
Interaction of G_i and cohort trends	0.483 (0.281)*	0.514 (0.277)*
No principal components	0.390 (0.201)*	0.444 (0.238)*
No never-takers	0.203 (0.253)	0.861 (0.243)***
Different outcome: delayed recall	0.115 (0.116)	0.353 (0.118)***
200 bootstrap repetitions	0.463 (0.219)**	0.471 (0.224)**

Notes: This table presents robustness checks for the linearized gene-environment estimates (our main result) using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. For reference, we provide our main result from Table 4. Robustness checks include calculating our main estimate (Eq. 8) over larger ranges of U_i^E , using only the most recent panel observation of each individual, relaxing the age restriction by keeping individuals below age 65, adding squared cohort trends, interacting cohort trends with G_i , excluding principal components (and their interaction with the instrument) from the control variables, excluding never-takers, and using the delayed word recall score (see Section 2.2) as an alternative outcome variable. Unless otherwise specified, standard errors are bootstrapped with 100 repetitions. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Nevertheless, we include robustness checks that add flexibility. First, we additionally interact cohort trends (including their interaction with the instrument) with G_i to allow trends to vary by PGI quintile. Doing so increases the linearized 2SLS coefficient of the $G \times E$ interaction marginally and raises the standard error. Second, we estimate a fully interacted model in which the controls are the levels and all interactions among cohort trends, the instrument, PGI quintiles, and the remaining previous controls (gender, principal components of genetic data, parental education, and survey wave fixed effects). The linearized estimate of this flexible model is somewhat closer to zero and has larger standard errors than our main result. Additionally, we estimate the 2SLS model with lower bandwidths of 8 and 5 years, respectively. The point estimate of the interaction changes marginally, and standard errors increase as the bandwidth decreases.

Finally, we show the robustness of our estimation method for the underlying optimization of the interaction effect (see Eq. 8). In our main specification, we maximize/minimize the difference between the first and the fifth quintile. Here, we also estimate the linearized effects of the final interaction when the optimization maximizes/minimizes the difference between the first and each of the other quintiles. We visualize the result in Figure D.3 in the Appendix. The results show that optimizing for different comparisons does not produce substantial changes in the final interaction effect, especially not in the crucial U_i^E range where we estimate our MTEs. The choice of quintile comparisons for the underlying optimization is not critical. The reason is that, irrespective of the specific quintile comparisons, the MTRs must match the same average outcomes across all response types and obey shape constraints, particularly additive separability, which identifies effects for the other quintiles that are unconsidered in the target parameter.

6 Conclusion

The growing gene-environment literature aims to estimate interactions between genetic endowments and environmental exposure (e.g., behaviour or choice variables like education) in their effect on an outcome of interest. The goal is to assess whether the effect of the environment varies by genetic endowment (or vice versa) while all else is equal. Since environmental variables are often endogenous, a popular choice is using instruments or natural experiments as a source of exogenous variation. This usually involves estimating a two-stage least squares model. Estimating gene-environment interactions by two-stage least squares regression identifies gene-specific effects of the environment. However, they may not retain

the desired interpretation as interaction effects if (1) the first stage is heterogeneous across different values of G_i and (2) the empirical setting entails essential heterogeneity in E_i (the unobserved heterogeneities for the outcome and treatment correlate). If both conditions hold, then two properties differ between gene-specific local average treatment effects: the genetic endowment and the unobserved effect heterogeneity. While the former is precisely what researchers want to isolate (the interaction), 2SLS cannot account for the latter. Thus, 2SLS estimates may not reflect complementarity between genes and the environment. We suggest solving this problem by estimating marginal treatment effects. MTEs allow for the computation of $G_i \times E_i$ estimates while accounting for unobserved heterogeneity.

While gene-environment interactions are a natural choice to illustrate this problem, since the central parameter is the instrumented interaction estimate, it theoretically applies to all interactions estimated by 2SLS. The two conditions that generate it, non-overlapping complier groups due to variations in the interaction variable and unobserved effect heterogeneity correlated with treatment propensity, could be present in other real-world scenarios involving choice variables. Nevertheless, there are likely also many settings where they are not present or the 2SLS comparisons are inconsequential. For example, [Barcellos *et al.* \(2021\)](#) find no differences between 2SLS and linear MTE estimates of their gene-education interaction. Moreover, in many applications, heterogeneous first-stage interactions with the interaction variable are unlikely, and studies that estimate only reduced-form (gene-environment) interaction effects avoid incorrect 2SLS comparisons altogether.

Our empirical study examines the long-term effects of education and genetic predisposition for education, as well as their interaction, on memory, our measure of cognition, using data from the English Longitudinal Study of Ageing. Word recall is frequently used as a

measure of cognitive functioning and predicts cognitive decline and impairment. To identify the effect of education, we use a compulsory schooling reform enacted in 1947 that raised the minimum school-leaving age in the UK to 15. Our baseline 2SLS estimates indicate no effect of education on recalled words among individuals in the lowest PGI quintile. Effects for higher quintiles are positive, but we lack the precision to estimate them precisely with 2SLS. We find evidence that both conditions for 2SLS to make the wrong comparisons apply in our setting. We see a strong gradient in the first stage across the quintiles of the education PGI and essential heterogeneity is present, more precisely, selection into gains. This is well-documented for educational decisions. We estimate marginal treatment effects using the partial identification approach from [Mogstad *et al.* \(2018\)](#). Building on reduced-form evidence, we generate minimal and maximal $G_i \times E_i$ effects consistent with the data. We add further benign restrictions (such as additive separability and negative MTE slopes that imply selection into gains) to gain precision and tighten the bounds. The resulting bounds almost point-identify the interaction effect.

Our main finding is that, holding unobserved heterogeneity across G_i fixed, even the lower-bound $G_i \times E_i$ effect is 4.7 times larger than the corresponding 2SLS estimate. In absolute terms, the gene-environment complementarity is substantial: on average, the effect of education on recalled words increases by 0.46–0.47 with each PGI quintile. This means that the MTE results imply higher returns to education for cognitive functioning later in life for those with a higher PGI. The complementarity between education and genetic predisposition widens existing gaps in returns to education – a finding that would remain undetected in our sample if estimated by two-stage least squares. This underscores that failing to account

for essential heterogeneity limits the usefulness of 2SLS estimates in our application—and beyond.

Acknowledgments

Financial support from Deutsche Forschungsgemeinschaft (DFG, project number 437564156) is gratefully acknowledged. This paper uses data from the English Longitudinal Study of Ageing (ELSA). ELSA is funded by the National Institute on Aging (R01AG017644), and by UK Government Departments coordinated by the National Institute for Health and Care Research (NIHR).

Affiliations

¹RWI - Leibniz Institute for Economic Research, Hohenzollernstr. 1–3, 45128 Essen,
Germany

²Paderborn University, Warburger Str. 100, 33098 Paderborn, Germany

³FernUniversität in Hagen, Universitätsstraße 47, 58097 Hagen, Germany

References

- Agostinelli, F. and Wiswall, M. (2025). ‘Estimating the technology of children’s skill formation’, *Journal of Political Economy*, vol. 133(3), pp. 846–887.
- Ahlskog, R., Beauchamp, J., Okbay, A., Oskarsson, S. and Thom, K. (2024). ‘Testing for treatment effect heterogeneity: Educational reform, genetic endowments, and family background’, Working paper, SSRN.
- Altonji, J.G. and Mansfield, R.K. (2018). ‘Estimating group effects using averages of observables to control for sorting on unobservables: School and neighborhood effects’, *American Economic Review*, vol. 108(10), pp. 2902–2946.
- Anderson, E.L., Howe, L.D., Wade, K.H., Ben-Shlomo, Y., Hill, W.D., Deary, I.J., Sanderson, E.C., Zheng, J., Korologou-Linden, R., Stergiakouli, E., Davey Smith, G., Davies, N.M. and Hemani, G. (2020). ‘Education, intelligence and Alzheimer’s disease: Evidence from a multivariable two-sample Mendelian randomization study’, *International Journal of Epidemiology*, vol. 49(4), pp. 1163–1172.
- Angrist, J.D. and Evans, W.N. (1998). ‘Children and their parents’ labor supply: Evidence from exogenous variation in family size’, *American Economic Review*, vol. 88(3), pp. 450–477.
- Angrist, J.D. and Krueger, A.B. (1991). ‘Does compulsory school attendance affect schooling and earnings?’, *The Quarterly Journal of Economics*, vol. 106(4), pp. 979–1014.
- Apolinario, D., Lichtenthaler, D.G., Magaldi, R.M., Soares, A.T., Busse, A.L., das Graças Amaral, J.R., Jacob-Filho, W. and Brucki, S.M.D. (2016). ‘Using temporal orientation,

category fluency, and word recall for detecting cognitive impairment: The 10-point cognitive screener (10-CS)', *International Journal of Geriatric Psychiatry*, vol. 31(1), pp. 4–12.

Arold, B.W., Hufe, P. and Stoeckli, M. (2025). 'Genetic endowments, educational outcomes and the mediating influence of school investments', *Journal of Political Economy Microeconomics*, forthcoming.

Banks, J., Batty, G.D., Breedvelt, J., Coughlin, K., Crawford, R., Marmot, M., Nazroo, J., Oldfield, Z., Steel, N., Steptoe, A., Wood, M. and Zaninotto, P. (2023). 'English Longitudinal Study of Ageing: Waves 0-9, 1998-2019', Data collection. 39th Edition. UK Data Service. SN: 5050.

Banks, J. and Mazzonna, F. (2012). 'The effect of education on old age cognitive abilities: Evidence from a regression discontinuity design', *The Economic Journal*, vol. 122(560), pp. 418–448.

Barcellos, S.H., Carvalho, L., Langa, K., Nimmagadda, S. and Turley, P. (2025). 'Education and dementia risk', Working paper 33430, National Bureau of Economic Research.

Barcellos, S.H., Carvalho, L. and Turley, P. (2021). 'The effect of education on the relationship between genetics, early-life disadvantages, and later-life ses', Working paper 28750, National Bureau of Economic Research.

Barcellos, S.H., Carvalho, L.S. and Turley, P. (2018). 'Education can reduce health differences related to genetic risk of obesity', *Proceedings of the National Academy of Sciences*, vol. 115(42), pp. E9765–E9772.

- Barth, D., Papageorge, N.W. and Thom, K. (2020). ‘Genetic endowments and wealth inequality’, *Journal of Political Economy*, vol. 128(4), pp. 1474–1522.
- Barth, D., Papageorge, N.W., Thom, K. and Velásquez-Giraldo, M. (2022). ‘Genetic endowments, income dynamics, and wealth accumulation over the lifecycle’, Working paper 30350, National Bureau of Economic Research.
- Behrman, J.R. and Taubman, P. (1989). ‘Is schooling ”mostly in the genes”? Nature–nurture decomposition using data on relatives’, *Journal of Political Economy*, vol. 97(6), pp. 1425–1446.
- Belsky, D.W., Domingue, B.W., Wedow, R., Arseneault, L., Boardman, J.D., Caspi, A., Conley, D., Fletcher, J.M., Freese, J., Herd, P., Moffitt, T.E., Poulton, R., Sicinski, K., Wertz, J. and Harris, K.M. (2018). ‘Genetic analysis of social-class mobility in five longitudinal studies’, *Proceedings of the National Academy of Sciences*, vol. 115(31), pp. E7275–E7284.
- Biroli, P., Galama, T.J., Hinke, S.v., Kippersluis, H.v., Rietveld, C.A. and Thom, K. (2025). ‘The economics and econometrics of gene–environment interplay’, *The Review of Economic Studies*, rdaf034.
- Björklund, A. and Salvanes, K.G. (2011). ‘Education and family background: Mechanisms and policies’, in (E. A. Hanushek, S. Machin and L. Woessmann, eds.), *Handbook of the Economics of Education*, pp. 201–247, vol. 3, Elsevier.
- Blundell, R. and Powell, J.L. (2003). ‘Endogeneity in nonparametric and semiparametric regression models’, in (M. Dewatripont, L. P. Hansen and S. J. Turnovsky, eds.), *Advances in economics and econometrics: Theory and applications, Eighth World Congress*, pp. 312–357, Econometric Society Monographs, Cambridge University Press.

- Bonsang, E., Adam, S. and Perelman, S. (2012). 'Does retirement affect cognitive functioning?', *Journal of Health Economics*, vol. 31(3), pp. 490–501.
- Brinch, C.N., Mogstad, M. and Wiswall, M. (2017). 'Beyond LATE with a discrete instrument', *Journal of Political Economy*, vol. 125(4), pp. 985–1039.
- Brunello, G., Weber, G. and Weiss, C.T. (2017). 'Books are forever: Early life conditions, education and lifetime earnings in europe', *The Economic Journal*, vol. 127(600), pp. 271–296.
- Bruno, D., Reiss, P.T., Petkova, E., Sidtis, J.J. and Pomara, N. (2013). 'Decreased recall of primacy words predicts cognitive decline', *Archives of Clinical Neuropsychology*, vol. 28(2), pp. 95–103.
- Cadar, D., Abell, J., Matthews, F.E., Brayne, C., Batty, G.D., Llewellyn, D.J. and Steptoe, A. (2020). 'Cohort profile update: The harmonised cognitive assessment protocol sub-study of the english longitudinal study of ageing (ELSA-HCAP)', *International Journal of Epidemiology*, vol. 50(3), pp. 725–726i.
- Card, D. (1993). 'Using geographic variation in college proximity to estimate the return to schooling', Working paper 4483, National Bureau of Economic Research.
- Carneiro, P., Heckman, J.J. and Vytlačil, E.J. (2011). 'Estimating marginal returns to education', *American Economic Review*, vol. 101(6), pp. 2754–2781.
- Carneiro, P. and Lee, S. (2009). 'Estimating distributions of potential outcomes using local instrumental variables with an application to changes in college enrollment and wage inequality', *Journal of Econometrics*, vol. 149(2), pp. 191–208.

- Carvalho, L.S. (2025). ‘Genetics and socioeconomic status: Some preliminary evidence on mechanisms’, *Journal of Political Economy Microeconomics*, vol. 3(3), pp. 429–476.
- Christelis, D., Jappelli, T. and Padula, M. (2010). ‘Cognitive abilities and portfolio choice’, *European Economic Review*, vol. 54(1), pp. 18–38.
- Clark, D. (2023). ‘School quality and the return to schooling in Britain: New evidence from a large-scale compulsory schooling reform’, *Journal of Public Economics*, vol. 223, p. 104902.
- Clark, D. and Royer, H. (2013). ‘The effect of education on adult mortality and health: Evidence from britain’, *American Economic Review*, vol. 103(6), pp. 2087–2120.
- Conti, G., Heckman, J. and Urzua, S. (2010). ‘The education–health gradient’, *American Economic Review*, vol. 100(2), pp. 234–238.
- Cornelissen, T., Dustmann, C., Raute, A. and Schönberg, U. (2018). ‘Who benefits from universal child care? Estimating marginal returns to early child care attendance’, *Journal of Political Economy*, vol. 126(6), pp. 2356–2409.
- Cunha, F. and Heckman, J.J. (2007). ‘The technology of skill formation’, *American Economic Review*, vol. 97(2), pp. 31–47.
- Currie, J. (2009). ‘Healthy, wealthy, and wise: Socioeconomic status, poor health in childhood, and human capital development’, *Journal of Economic Literature*, vol. 47(1), pp. 87–122.
- Devereux, P.J. and Hart, R.A. (2010). ‘Forced to be rich? Returns to compulsory schooling in britain’, *The Economic Journal*, vol. 120(549), pp. 1345–1364.

- Ding, X., Barban, N., Tropf, F.C. and Mills, M.C. (2019). ‘The relationship between cognitive decline and a genetic predictor of educational attainment’, *Social Science & Medicine*, vol. 239, p. 112549.
- Gathmann, C., Jürges, H. and Reinhold, S. (2015). ‘Compulsory schooling reforms, education and mortality in twentieth century Europe’, *Social Science & Medicine*, vol. 127, pp. 74–82.
- Glymour, M.M., Kawachi, I., Jencks, C.S. and Berkman, L.F. (2008). ‘Does childhood schooling affect old age memory or mental status? Using state schooling laws as natural experiments’, *Journal of Epidemiology & Community Health*, vol. 62(6), pp. 532–537.
- Gorman, E. (2023). ‘Does schooling have lasting effects on cognitive function? Evidence from compulsory schooling laws’, *Demography*, vol. 60(4), pp. 1139–1161.
- Gottfredson, L.S. (1997). ‘Mainstream science on intelligence: An editorial with 52 signatories, history, and bibliography’, *Intelligence*, vol. 24(1), pp. 13–23.
- Harmon, C. and Walker, I. (1995). ‘Estimates of the economic return to schooling for the United Kingdom’, *American Economic Review*, vol. 85(5), pp. 1278–1286.
- Heckman, J.J., Urzua, S. and Vytlačil, E. (2006). ‘Understanding instrumental variables in models with essential heterogeneity’, *The Review of Economics and Statistics*, vol. 88(3), pp. 389–432.
- Heckman, J.J. and Vytlačil, E. (2005). ‘Structural equations, treatment effects, and econometric policy evaluation’, *Econometrica*, vol. 73(3), pp. 669–738.

- Houmark, M.A., Ronda, V. and Rosholm, M. (2024). ‘The nurture of nature and the nature of nurture: How genes and investments interact in the formation of skills’, *American Economic Review*, vol. 114(2), pp. 385–425.
- Imbens, G.W. and Angrist, J.D. (1994). ‘Identification and estimation of local average treatment effects’, *Econometrica*, vol. 62(2), pp. 467–475.
- Imbens, G.W. and Manski, C.F. (2004). ‘Confidence intervals for partially identified parameters’, *Econometrica*, vol. 72(6), pp. 1845–1857.
- Imbens, G.W. and Newey, W.K. (2009). ‘Identification and estimation of triangular simultaneous equations models without additivity’, *Econometrica*, vol. 77(5), pp. 1481–1512.
- Imbens, G.W. and Rubin, D.B. (1997). ‘Estimating outcome distributions for compliers in instrumental variables models’, *The Review of Economic Studies*, vol. 64(4), pp. 555–574.
- Ito, K., Ida, T. and Tanaka, M. (2023). ‘Selection on welfare gains: Experimental evidence from electricity plan choice’, *American Economic Review*, vol. 113(11), pp. 2937–2973.
- Johnston, D.W., Lordan, G., Shields, M.A. and Suziedelyte, A. (2015). ‘Education and health knowledge: Evidence from UK compulsory schooling reform’, *Social Science & Medicine*, vol. 127, pp. 92–100.
- Jürges, H., Kruk, E. and Reinhold, S. (2013). ‘The effect of compulsory schooling on health—evidence from biomarkers’, *Journal of Population Economics*, vol. 26(2), pp. 645–672.

- Kamhöfer, D.A., Schmitz, H. and Westphal, M. (2019). ‘Heterogeneity in marginal non-monetary returns to higher education’, *Journal of the European Economic Association*, vol. 17(1), pp. 205–244.
- Kline, P. and Walters, C.R. (2019). ‘On heckits, LATE, and numerical equivalence’, *Econometrica*, vol. 87(2), pp. 677–696.
- Kowalski, A.E. (2023). ‘Reconciling seemingly contradictory results from the Oregon Health Insurance Experiment and the Massachusetts Health Reform’, *The Review of Economics and Statistics*, vol. 105(3), pp. 646–664.
- Krumme, A. and Westphal, M. (2024). ‘Monetary returns to upper secondary schooling, the evolution of unobserved heterogeneity, and implications for employer learning’, Ruhr Economic Papers #1130, RWI - Leibniz Institute for Economic Research.
- Lee, J.J., Wedow, R., Okbay, A., Kong, E., Maghzian, O., Zacher, M., Nguyen-Viet, T.A., Bowers, P., Sidorenko, J., Karlsson Linnér, R., Fontana, M.A., Kundu, T., Lee, C., Li, H., Li, R., Royer, R., Timshel, P.N., Walters, R.K., Willoughby, E.A., Yengo, L., Alver, M., Bao, Y., Clark, D.W., Day, F.R., Furlotte, N.A., Joshi, P.K., Kemper, K.E., Kleinman, A., Langenberg, C., Mägi, R., Trampush, J.W., Verma, S.S., Wu, Y., Lam, M., Zhao, J.H., Zheng, Z., Boardman, J.D., Campbell, H., Freese, J., Harris, K.M., Hayward, C., Herd, P., Kumari, M., Lencz, T., Luan, J., Malhotra, A.K., Metspalu, A., Milani, L., Ong, K.K., Perry, J.R.B., Porteous, D.J., Ritchie, M.D., Smart, M.C., Smith, B.H., Tung, J.Y., Wareham, N.J., Wilson, J.F., Beauchamp, J.P., Conley, D.C., Esko, T., Lehrer, S.F., Magnusson, P.K.E., Oskarsson, S., Pers, T.H., Robinson, M.R., Thom, K., Watson, C., Chabris, C.F., Meyer, M.N., Laibson, D.I., Yang, J., Johannesson, M., Koellinger, P.D.,

- Turley, P., Visscher, P.M., Benjamin, D.J. and Cesarini, D. (2018). ‘Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals’, *Nature Genetics*, vol. 50(8), pp. 1112–1121.
- Maestas, N., Mullen, K.J. and Strand, A. (2013). ‘Does disability insurance receipt discourage work? Using examiner assignment to estimate causal effects of SSDI receipt’, *American Economic Review*, vol. 103(5), pp. 1797–1829.
- Manski, C.F. (1997). ‘Monotone treatment response’, *Econometrica*, vol. 65(6), pp. 1311–1334.
- Mogstad, M., Santos, A. and Torgovitsky, A. (2018). ‘Using instrumental variables for inference about policy relevant treatment parameters’, *Econometrica*, vol. 86(5), pp. 1589–1619.
- Mogstad, M. and Torgovitsky, A. (2018). ‘Identification and extrapolation of causal effects with instrumental variables’, *Annual Review of Economics*, vol. 10(1), pp. 577–613.
- Mogstad, M. and Torgovitsky, A. (2024). ‘Instrumental variables with unobserved heterogeneity in treatment effects’, in (C. Dustmann and T. Lemieux, eds.), *Handbook of Labor Economics*, pp. 1–114, vol. 5, Elsevier.
- Muslimova, D., Van Kippersluis, H., Rietveld, C.A., Von Hinke, S. and Meddens, S.F.W. (2025). ‘Gene–environment complementarity in educational attainment’, *Journal of Labor Economics*, forthcoming.
- Nybom, M. (2017). ‘The distribution of lifetime earnings returns to college’, *Journal of Labor Economics*, vol. 35(4), pp. 903–952.

- Oreopoulos, P. (2006). 'Estimating average and local average treatment effects of education when compulsory schooling laws really matter', *American Economic Review*, vol. 96(1), pp. 152–175.
- Papageorge, N.W. and Thom, K. (2020). 'Genes, education, and labor market outcomes: Evidence from the Health and Retirement Study', *Journal of the European Economic Association*, vol. 18(3), pp. 1351–1399.
- Pereira, R.D., Rietveld, C.A. and van Kippersluis, H. (2025). 'The interplay between maternal smoking and genes in offspring birth weight', *Journal of Human Resources*, vol. 60(2), pp. 400–433.
- Plomin, R. (2014). 'Genotype–environment correlation in the Era of DNA', *Behavior Genetics*, vol. 44(6), pp. 629–638.
- Plomin, R., DeFries, J.C. and Loehlin, J.C. (1977). 'Genotype-environment interaction and correlation in the analysis of human behavior', *Psychological Bulletin*, vol. 84(2), pp. 309–322.
- Plug, E. and Vijverberg, W. (2003). 'Schooling, family background, and adoption: Is it nature or is it nurture?', *Journal of Political Economy*, vol. 111(3), pp. 611–641.
- Powdthavee, N. (2010). 'Does education reduce the risk of hypertension? Estimating the biomarker effect of compulsory schooling in England', *Journal of Human Capital*, vol. 4(2), pp. 173–202.

- Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A. and Reich, D. (2006). 'Principal components analysis corrects for stratification in genome-wide association studies', *Nature Genetics*, vol. 38(8), pp. 904–909.
- Rohwedder, S. and Willis, R.J. (2010). 'Mental retirement', *Journal of Economic Perspectives*, vol. 24(1), pp. 119–138.
- Rose, E.K. and Shem-Tov, Y. (2021). 'How does incarceration affect reoffending? Estimating the dose–response function', *Journal of Political Economy*, vol. 129(12), pp. 3302–3356.
- Roy, A.D. (1951). 'Some thoughts on the distribution of earnings', *Oxford Economic Papers*, vol. 3(2), pp. 135–146.
- Schiele, V. and Schmitz, H. (2023). 'Understanding cognitive decline in older ages: The role of health shocks', *European Economic Review*, vol. 151, p. 104320.
- Schmitz, H. and Westphal, M. (2025). 'Early- and later-life stimulation: how retirement shapes the effect of education on old-age cognitive abilities', Ruhr Economic Papers #1146, RWI - Leibniz Institute for Economic Research.
- Schmitz, L.L. and Conley, D. (2017). 'The effect of Vietnam-era conscription and genetic potential for educational attainment on schooling outcomes', *Economics of Education Review*, vol. 61, pp. 85–97.
- Silles, M.A. (2009). 'The causal effect of education on health: Evidence from the United Kingdom', *Economics of Education Review*, vol. 28(1), pp. 122–128.

Steptoe, A., Breeze, E., Banks, J. and Nazroo, J. (2013). ‘Cohort profile: The English Longitudinal Study of Ageing’, *International Journal of Epidemiology*, vol. 42(6), pp. 1640–1648.

Tsoi, K.K.F., Chan, J.Y.C., Hirai, H.W., Wong, A., Mok, V.C.T., Lam, L.C.W., Kwok, T.C.Y. and Wong, S.Y.S. (2017). ‘Recall tests are effective to detect mild cognitive impairment: A systematic review and meta-analysis of 108 diagnostic studies’, *Journal of the American Medical Directors Association*, vol. 18(9), pp. 807.e17–807.e29.

van den Berg, G.J., von Hinke, S. and Vitt, N. (2023a). ‘Early life exposure to measles and later-life outcomes: Evidence from the introduction of a vaccine’, ArXiv preprint, arXiv:2301.10558.

van den Berg, G.J., von Hinke, S. and Wang, R.A.H. (2023b). ‘Prenatal sugar consumption and late-life human capital and health: Analyses based on Postwar rationing and polygenic scores’, ArXiv preprint, arXiv:2301.09982.

Westphal, M., Kamhöfer, D.A. and Schmitz, H. (2022). ‘Marginal college wage premiums under selection into employment’, *The Economic Journal*, vol. 132(646), pp. 2231–2272.

A Implications for other empirical applications

Of course, the problem described in this paper and its solution are not limited to gene-environment interactions. Many applied papers focus on observed heterogeneity in treatment effects in general and on the effect gradient along one particular variable. Consider the following regression equation:

$$Y_i = \beta_0 + \beta_1 D_i + \beta_2 X_i + \beta_3 D_i \times X_i + \varepsilon_i, \quad (\text{A.1})$$

where Y_i is the outcome, D_i is the (endogenous) treatment indicator, and X_i is an observable variable. We briefly mention several studies and fields in the literature that are interested in how the effect of treatment D varies by X . However, we do not claim that problems with 2SLS are always present because the conditions under which 2SLS causes problems are not always met.

Education

The literature on effects of education D has been interested in how they vary by socio-economic status X . For instance, early research on local average treatment effects of education focused on identifying effect gradients by socio-economic status to provide additional support for the plausibility of the IV estimates ([Angrist and Krueger, 1991](#); [Card, 1993](#)). Relatedly, [Brunello *et al.* \(2017\)](#) want to identify the observed heterogeneity in the returns to education explicitly. They distinguish between urban and rural areas, as well as between individuals who grow up in households with many books versus those with few. They interpret their heterogeneous returns to compulsory schooling both structurally and

in terms of selection (via different marginal costs and returns of schooling for affected individuals). Nonetheless, applying our method in this context could help differentiate between the two explanations. Furthermore, to gain better insights into possible structural disadvantages of women, researchers interested in the effects of education by gender should compare average treatment effects across genders rather than local average treatment effects for two possibly different complier groups.

Children

Relatedly, D could be the effect of children, Y labour supply, X household income: In their landmark study, [Angrist and Evans \(1998\)](#) use the same-sex instrument to estimate whether the effects of children on labour supply differ by household income (as economic theory would predict). To assess the validity of this theory, one would need to hold the complier population constant and compare the effects of children on labour supply between different groups formed by household income. Given that [Angrist and Evans \(1998\)](#) show that the first stage is increasing along the husband's income (the more affluent the household, the more the household can afford their same-sex preference), 2SLS might be problematic here.

Skill formation

Other examples relate to the study of origins of the socio-economic gradient in health ([Currie, 2009](#)) and the literature on skill formation ([Cunha and Heckman, 2007](#)). These are prominent examples of a research question where the interaction effect between the environment in which individuals grow up (X) and a choice variable D is at the centre of interest. For instance, [Conti et al. \(2010\)](#) estimate the interaction between education and measures (X) of early-life cognitive and non-cognitive skills, as well as health endowments (the dimensions of observed heterogeneity) on later-life health behaviour using structural methods (such as

latent factor models that allow proxying unobservables and measurement error). If one were to employ instrumental variables estimation for this research question, it would be essential to consider differential responses to the instruments for all observed heterogeneity dimensions.

Another example of this literature is the recent contribution by [Agostinelli and Wiswall \(2025\)](#). They specify a latent factor model that explores the heterogeneity in returns to parental investments concerning children's endowments and propose an empirical specification. If quasi-experimental instruments for parental investments exist, and if endowments and investments can be observed directly without measurement error, then this model can be estimated using instrumental variable methods. However, the problems outlined in this paper would arise if there were a gradient in the first stage based on children's endowments, or if parents, for instance, had a higher propensity to invest more in children with high returns on their investment.

Related to our research question, [Houmark *et al.* \(2024\)](#) estimates the technology of skill formation using skills, genetic endowment of children and parents, and parental investments, documenting that all factors are interrelated.

Social security benefits

Heterogeneity between unobserved heterogeneity and covariates also appears to interfere in the paper by [Maestas *et al.* \(2013\)](#). They estimate the effects of disability benefit receipt on employment and report heterogeneous first-stage and IV regressions that vary by observed characteristics. For the first stage, coefficients on the instrument and the intercept vary substantially by covariates, suggesting that the instrument affects complier groups with different unobserved characteristics. While interpreting the 2SLS coefficients as heterogeneous

LATEs is (of course) appropriate (as they do), this suggests that generalizing these LATEs beyond the instrument-specific complying population, as covariate-specific average treatment effects or covariate-specific effects of a more lenient or stringent disability receipt allowance reform, is likely inappropriate. This is especially important because [Maestas *et al.* \(2013\)](#) detect essential heterogeneity: employment effects of disability receipt are less negative for individuals with a higher unobserved severity.

B Polygenic indices

The human genome has about 3 billion base pairs, the pairs of nucleic acids that make up the DNA. However, any two people differ by only about 0.1 % of the base pairs. Most of these genetic differences are substitutions of a single base (adenine, thymine, cytosine, or guanine) for another at a specific location in the genome, called "single nucleotide polymorphisms" (SNPs) that are common across the whole genome. These substitutions result in different genetic variants (alleles) that vary among parts of the population.¹⁷ For example, at a specific SNP location, the DNA sequence might have an adenine base in some individuals, while others may have a thymine base at the same position. One is (arbitrarily) chosen as the reference variant. Then, each SNP can be represented as a count variable of occurrences of the reference variant at this location that can either be 0, 1 or 2, since there are two copies of each chromosome. Large research projects called genome-wide association studies (GWAS) correlate each of the SNPs with a disease or trait, e.g., diabetes, years of education, or smoking. This entails running one regression of type

$$Y_i = \beta_j S_{ij} + X_i' \delta + \zeta_i \quad (\text{B.2})$$

for each of the SNPs, where Y_i is the outcome of interest (in our case educational attainment) of individual i , β_j is the individual effect of each SNP j , S_{ij} is the count variable of the reference variant of SNP j with $S_{ij} \in \{0, 1, 2\}$, X_i is a vector of controls that typically include age, gender and principal components of the genetic data, which control for spurious

¹⁷The generally agreed-upon threshold for a substitution to be regarded a SNP is common occurrence in at least one % of the population.

correlations of genetic variants and outcomes that are due to population structure.¹⁸ The PGI is then calculated as a weighted sum of all J SNPs that are relevant to the outcome¹⁹, where the weights correspond to the β_j 's obtained in the GWAS:

$$PGI_i = \sum_{j=1}^J \beta_j S_{ij} \quad (\text{B.3})$$

Polygenic scores for various traits or behaviours (personality, mental and physical health, health behaviours, and more) have been calculated for the ELSA sample based on various GWAS and are readily available.

¹⁸For a more detailed description of principal components, see Section 2.2.

¹⁹The discovery study the PGI we are using is based on, [Lee *et al.* \(2018\)](#), finds 1,271 SNPs that are significantly correlated with educational attainment.

C Interpreting $G \times Z$

In the $G \times E$ literature, many studies use an exogenous environment, such as a policy change, which we refer to as Z_i . Hence, these studies estimate reduced-form $G \times Z$ interactions. [Muslimova *et al.* \(2025\)](#), for instance, use a firstborn indicator as the measure of environment, which is not an individual decision, and hence, exogenous to the individual. In other studies, not every individual is affected by a change in Z_i . Examples include [Schmitz and Conley \(2017\)](#), where Z_i is the Vietnam draft lottery which provides incentives for education, the implicit E_i), [van den Berg *et al.* \(2023a\)](#), where Z_i is a vaccination campaign, and the implicit E_i would be measles infections), [van den Berg *et al.* \(2023b\)](#), with Z_i constituting a sugar derationing policy, with the implicit E_i being the maternal sugar consumption), and [Ahlskog *et al.* \(2024\)](#), where Z_i is a compulsory schooling reform shifting education, as in our setting).

As in our study, Z_i constitutes an incentive for the underlying individual decision (or behaviour) E_i we are interested in. While the focus on $G \times Z$ can be the policy-relevant effect (depending on the context), the focus on the reduced form does not solve problems with essential heterogeneity. As shown now, the $G \times Z$ interaction may be driven solely by a first-stage gradient, even if the (latent structural) $G \times E$ interaction is absent.

Rearranging the Wald estimator demonstrates that the reduced-form effect is the product of the structural (i.e., the second stage) and the first stage effects. This expression holds at every conceivable value of $G = g$ (which, for the sake of the argument, is assumed to be continuous):

$$\begin{aligned} \mathbb{E}(Y \mid Z = 1, G = g) - \mathbb{E}(Y \mid Z = 0, G = g) &= \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g)) \\ &\quad \times \mathbb{E}(\mathbf{1}[C(G = g)]) \end{aligned} \tag{C.4}$$

Conceptually, you can think of this equation as the following form: $w(g) := u(g, v(g)) \times v(g)$, where we define the reduced form $w(g)$ as a function of g . It is the product of $u(\cdot)$ —the second stage—and $v(g)$ —the first stage (if aggregated). The second stage varies in $G = g$, the PGI at which the effect is assessed, and the specific complier group determined by $v(\cdot)$, a complier-indicating function (if unaggregated).

Applying the product and chain rule to this simplified expression $w(g)$ demonstrates that marginally changing g has three distinct effects (i.e., partial derivatives) on the reduced form: (i) the partial direct derivative of $w(\cdot)$ with respect to the structural gene heterogeneity of the second stage (i.e., $\frac{\partial u(\cdot)}{\partial g}$, holding $v(g)$ fixed), (ii) the partial derivative of $u(g)$ with respect to the (g -specific) complier groups (i.e. $\frac{\partial u(\cdot)}{\partial v(\cdot)}$ holding g fixed), and (iii) the partial derivative of g with respect to the first stage $\frac{\partial v(\cdot)}{\partial g}$.

Formally in the notation of Eq. (C.4), this reads:

$$\begin{aligned}
& \frac{\partial \mathbb{E}(Y \mid Z = 1, G = g) - \mathbb{E}(Y \mid Z = 0, G = g)}{\partial g} = \\
& \underbrace{\frac{\partial \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g))}{\partial g}}_{\text{(i) Structural outcome interaction}} \times \mathbb{E}\left(\mathbb{1}[C(G = g)]\right) \\
& + \underbrace{\frac{\partial \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g))}{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}}_{\text{(ii) Interference with essential heterogeneity}} \\
& \times \frac{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}{\partial g} \times \mathbb{E}(C(G = g)) \\
& + \underbrace{\mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g))}_{\text{(iii) First-stage gradient}} \times \overbrace{\frac{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}{\partial g}}^{<0 \text{ in our paper}}
\end{aligned}$$

The interaction between channels (i) and (ii) is the core of our paper. The focus on the reduced-form adds a third channel, which blurs the structural outcome interaction we want to identify. To see this clearly, assume there is no structural outcome interaction and also no essential heterogeneity—so channels (i) and (ii) are switched off. Yet, the first stage may exhibit a gradient. In this case, the reduced form interaction $G \times Z$ may still differ from zero. It remains:

$$\begin{aligned}
& \frac{\partial \mathbb{E}(Y \mid Z = 1, G = g) - \mathbb{E}(Y \mid Z = 0, G = g)}{\partial g} = \mathbb{E}(Y^1(G = g) - Y^0(G = g) \mid C(G = g)) \\
& \times \frac{\partial \mathbb{E}\left(\mathbb{1}[C(G = g)]\right)}{\partial g}
\end{aligned}$$

This demonstrates that even without any structural outcome interaction and essential heterogeneity, the $G \times Z$ interaction may differ from zero. This is because the first stage changes along G . While the $G \times Z$ may be informative about whether a policy Z has heterogeneous effects along G (as measured by the sum of all three mechanisms), it is uninformative about the structural outcome interaction (mechanism (i)).

D Additional Tables and Figures

Table D.1: Descriptive statistics by availability of genetic information

	Full sample	By availability of genetic information		
	Mean (SD)	Yes	No	Difference (SE)
<i>Outcome Y_i</i>				
Word recall score	9.32 (3.50)	9.67	8.76	0.91 (0.05)***
<i>Treatment E_i</i>				
Left school ≥ 15	0.76 (0.43)	0.78	0.72	0.06 (0.01)***
<i>Instrument Z_i</i>				
Born 1933 or later	0.65 (0.48)	0.66	0.62	0.04 (0.01)***
<i>Controls</i>				
Female	0.52 (0.50)	0.52	0.52	0.00 (0.01)
Birth year	1934.67 (5.10)	1934.89	1934.32	0.57 (0.08)***
Parental education:				
Missing	0.31 (0.46)	0.25	0.40	-0.16 (0.01)***
Both left school ≤ 14	0.53 (0.50)	0.57	0.45	-0.12 (0.01)***
At least one left school ≥ 15	0.17 (0.37)	0.18	0.14	0.04 (0.01)***
Wave (years):				
- 1 (2002-2003)	0.19 (0.39)	0.15	0.24	-0.09 (0.01)***
- 2 (2004-2005)	0.17 (0.38)	0.18	0.15	0.03 (0.01)***
- 3 (2006-2007)	0.15 (0.36)	0.17	0.13	0.04 (0.01)***
- 4 (2008-2009)	0.19 (0.39)	0.19	0.18	0.01 (0.01)
- 5 (2010-2011)	0.16 (0.37)	0.17	0.16	0.01 (0.01)
- 6 (2012-2013)	0.14 (0.34)	0.14	0.13	0.01 (0.01)*
Observations	17,884	11,027	6,857	

Notes: This table presents descriptive statistics by availability of genetic information in ELSA using data from waves 1–6 and different sample restrictions. “Full sample” includes all restrictions outlined in Section 2.2, except for the removal of individuals without genetic information. In the second part of this table, we split this sample based on the availability of genetic information. The sub-sample with genetic information corresponds to our main estimation sample. We display the means and difference of means, as well as the standard errors of a t-test for equality of means between the two groups. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Table D.2: The 1947 UK compulsory schooling reform and providing genetic information to ELSA

	DV: Provided genetic information (1)	DV: Left school at 15 or later (E_i) (2)
Z_i	-0.01 (0.030)	0.479 (0.030)***
Provided genetic information		0.035 (0.023)
Provided genetic information $\times Z_i$		-0.008 (0.025)
Controls	Yes	Yes
Observations	17,884	17,884

Notes: This table shows that our instrument, eligibility for the 1947 UK compulsory schooling reform, did not affect the probability of providing genetic information to ELSA (column 1) and that the first stage does not vary with the provision of genetic information to ELSA (column 2) using data from ELSA waves 1–6 and the sample selection outlined in Section 2.2, except for the removal of individuals without genetic information. Specifically, column 1 shows the estimates of a linear regression of the instrument Z_i on a dummy variable equal to one if a person provided genetic information to ELSA, and column 2 shows the estimates of a linear regression of the environment (leaving school at age 15 or later) on the instrument Z_i , being born in 1933 or later, the genetic information dummy, and the interaction between the two. The controls in each case include a linear cohort trend, its interaction with the instrument, gender, and survey wave fixed effects. Both regressions are estimated using a larger sample from Table D.1, where we apply all the sample restrictions outlined in Section 2.2, except for removing individuals without genetic information. Standard errors in both regressions are clustered at the individual level. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Table D.3: Descriptive statistics (extended)

	Main sample	By E_i		
	Mean (SD)	$E_i=1$	$E_i=0$	Difference (SE)
<i>Outcome Y_i</i>				
Word recall score	9.67 (3.37)	10.11	8.08	2.03 (0.07)***
<i>Treatment E_i</i>				
Left school ≥ 15	0.78 (0.41)	1.00	0.00	1.00 (0.00)
<i>Polygenic index G_i</i>				
1st PGI quintile	0.20 (0.40)	0.18	0.25	-0.07 (0.01)***
2nd PGI quintile	0.19 (0.40)	0.19	0.21	-0.02 (0.01)**
3rd PGI quintile	0.20 (0.40)	0.21	0.19	0.02 (0.01)**
4th PGI quintile	0.21 (0.41)	0.21	0.20	0.01 (0.01)
5th PGI quintile	0.20 (0.40)	0.22	0.15	0.07 (0.01)***
<i>Instrument Z_i</i>				
Born 1933 or later	0.66 (0.47)	0.82	0.13	0.69 (0.01)***
<i>Controls</i>				
Female	0.52 (0.50)	0.52	0.50	0.02 (0.01)**
Birth year	1934.89 (5.00)	1936.29	1929.92	6.37 (0.10)***
Parental education:				
Missing	0.25 (0.43)	0.20	0.41	-0.21 (0.01)***
Both left school ≤ 14	0.57 (0.49)	0.58	0.55	0.03 (0.01)**
At least one left school ≥ 15	0.18 (0.39)	0.22	0.04	0.18 (0.01)***
Principal components (standardized):				
- 1 -	0.00 (1.00)	0.00	-0.01	0.02 (0.02)
- 2 -	0.00 (1.00)	0.01	-0.02	0.03 (0.02)
- 3 -	0.00 (1.00)	0.01	-0.04	0.05 (0.02)**
- 4 -	0.00 (1.00)	-0.01	0.02	-0.03 (0.02)
- 5 -	0.00 (1.00)	0.00	0.00	0.00 (0.02)
- 6 -	0.00 (1.00)	0.02	-0.07	0.09 (0.02)***
- 7 -	0.00 (1.00)	0.01	-0.03	0.04 (0.02)*
- 8 -	0.00 (1.00)	0.00	0.02	-0.02 (0.02)
- 9 -	0.00 (1.00)	0.01	-0.02	0.02 (0.02)
- 10 -	0.00 (1.00)	0.01	-0.02	0.02 (0.02)
Wave (years):				
- 1 (2002-2003)	0.15 (0.36)	0.12	0.27	-0.15 (0.01)***
- 2 (2004-2005)	0.18 (0.38)	0.16	0.26	-0.10 (0.01)***
- 3 (2006-2007)	0.17 (0.37)	0.16	0.19	-0.03 (0.01)***
- 4 (2008-2009)	0.19 (0.39)	0.20	0.15	0.06 (0.01)***
- 5 (2010-2011)	0.17 (0.37)	0.19	0.09	0.10 (0.01)***
- 6 (2012-2013)	0.14 (0.35)	0.17	0.04	0.13 (0.01)***
Observations	11,027	8,590	2,437	

Notes: This table presents extended descriptive statistics using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. Here, we also report the standardized first 10 principal components of the genetic data and information on survey waves. The categories for parental education include: Missing information of at least one parent, both parents left full-time education at age 14 or before or have no education, and at least one parent stayed in school until age 15 or longer. We include the mean and standard deviation of the main sample as well as the means by E_i , the difference of means, and the standard errors of a t-test for equality of means. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Table D.4: Descriptive statistics by PGI quintiles

	1st quintile	2nd quintile	3rd quintile	4th quintile	5th quintile
<i>Outcome Y_i</i>					
Word recall score	8.98	9.47	9.69	9.84	10.31
<i>Treatment E_i</i>					
Left school ≥ 15	0.72	0.76	0.80	0.78	0.84
<i>Instrument Z_i</i>					
Born 1933 or later	0.67	0.65	0.66	0.67	0.67
<i>Controls</i>					
Female	0.54	0.51	0.53	0.50	0.50
Birth year	1934.87	1934.84	1935.05	1934.69	1934.98
Parental education:					
Missing	0.27	0.26	0.22	0.23	0.25
Both left school ≤ 14	0.62	0.57	0.59	0.62	0.47
At least one left school ≥ 15	0.11	0.17	0.19	0.15	0.29
Wave (years):					
– 1 (2002-2003)	0.16	0.15	0.14	0.16	0.15
– 2 (2004-2005)	0.19	0.18	0.17	0.19	0.17
– 3 (2006-2007)	0.17	0.17	0.17	0.17	0.17
– 4 (2008-2009)	0.19	0.19	0.20	0.18	0.20
– 5 (2010-2011)	0.16	0.17	0.17	0.16	0.17
– 6 (2012-2013)	0.13	0.14	0.15	0.14	0.15
Observations	2,152	2,145	2,216	2,284	2,230

Notes: This table presents sample means by quintiles of the educational attainment PGI using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2.

Table D.5: The 1947 UK compulsory schooling reform and panel attrition

	DV: Dropped out of sample	
	Coefficient (1)	Standard error (2)
Z_i	0.001	(0.018)
Controls	Cohort trends only	
Observations	12,108	

Notes: This table presents estimates of the effect of the 1947 UK compulsory schooling reform (Z_i) on a panel attrition indicator. The analysis uses data from ELSA waves 1–6, our main sample selection, as outlined in Section 2.2, but we fill up observations from the first wave an individual was observed until wave 6 to create the panel attrition indicator. Controls include a linear cohort trend and its interaction with the instrument. Standard errors clustered at the individual level shown are in parentheses. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Table D.6: Estimates of the first stages by PGI quintile

	DV: Left school at 15 or later (E_i)	
	Coefficient (1)	Standard error (2)
$Z_i \times (G_i = 1)$	0.649	(0.018)***
$Z_i \times (G_i = 2)$	0.534	(0.018)***
$Z_i \times (G_i = 3)$	0.484	(0.018)***
$Z_i \times (G_i = 4)$	0.426	(0.018)***
$Z_i \times (G_i = 5)$	0.357	(0.018)***
Controls	Yes	
Observations	11,027	

Notes: This table presents estimates of the effect of the 1947 UK compulsory schooling reform (Z_i) on attending school until at least age 15 (E_i) by quintiles of the education PGI using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. These effects are obtained from the coefficients $\pi_{1,\Delta}^q$ to $\pi_{5,\Delta}^q$ of Eq. (4), which correspond to the complier shares in the respective quintile. Standard errors clustered at the individual level shown are in parentheses. The controls include a linear cohort trend, its interaction with the instrument, gender, survey wave fixed effects, parental education, the first ten principal components of the genetic data as well as their interactions with the instrument. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$. The standard errors are seemingly the same here, but only due to rounding. They differ at the fourth digit.

Table D.7: MTE estimation without interaction terms

	DV: Total word recall score	
	MTE _{min} (1)	MTE _{max} (2)
E_i	0.163 (0.464)	1.636 (0.620)***
Controls	Yes	Yes
Observations	11,027	11,027

Notes: This table shows MTE estimates of the total effect of E_i on word recall, i.e., without interacting E_i with the genetic endowment. The analysis uses data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. We use our method as described in Section 5, but only construct two MTE curves, that maximize/minimize the total effect instead of two for each PGI quintile. Controls include a linear cohort trend, its interaction with the instrument, gender, and survey wave fixed effects. Standard errors in both regressions are bootstrapped with 100 repetitions. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Table D.8: The gradient within essential heterogeneity of other studies

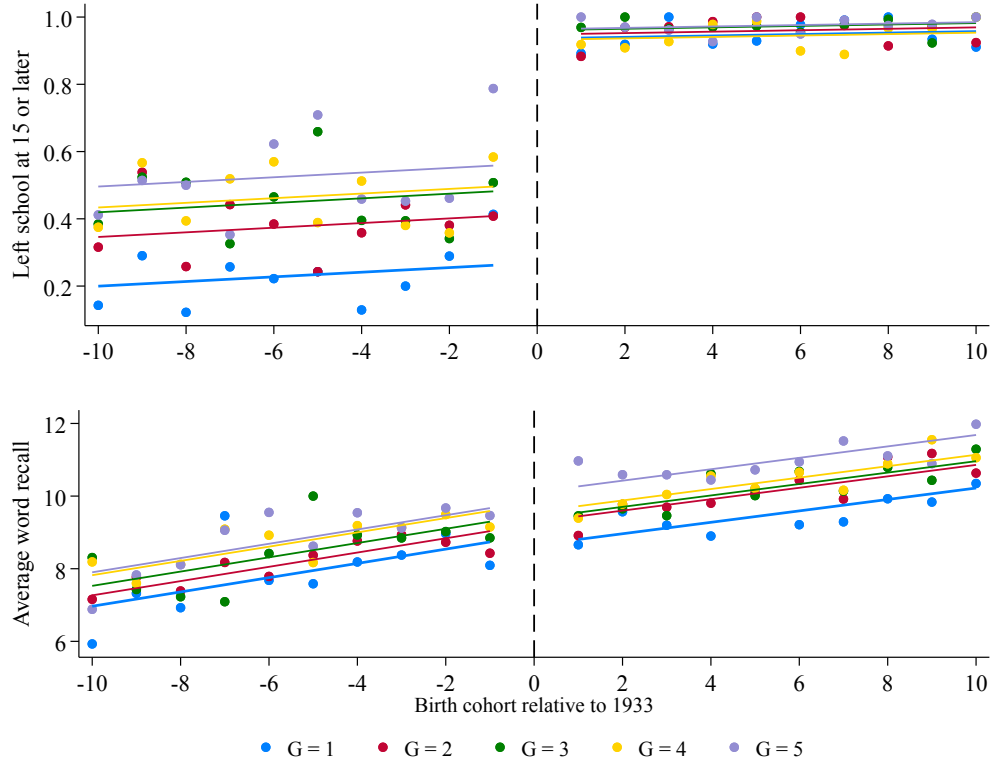
Study	Outcome	Treatment	Note (Model/sample)	(1) $MTE(U^E = p)$ $p = 0.1$	(2) $p = 0.9$	(3) Linearized effect per quintile $(\frac{(1)-(2)}{5})$	(4) Standard deviation of outcome	(5) Standardized linearized effect size $(\frac{(3)}{(4)})$	(6) % of our MSE size $(\frac{(5)}{(0.12)})$
Nybohm (2017)	Lifetime earnings	College ed.	Norm. sel. mod. Local IV	0.04 0.08	-0.02 0.1	0.015 0.005	0.4 0.4	0.03 0.0125	37.5% 10.4%
Carneiro et al. (2011)	Wage in 1991	College ed.	Norm. sel. mod. Local IV	0.14 0.35	0 -0.15	0.035 0.1	0.47 0.47	0.06 0.21	62% 175%
Kamhöfer et al. (2019)	Math literacy Reading speed Reading competence	College ed.	Local IV Local IV Local IV	2.00 1.8 2.5	0.5 0.3 0	0.375 0.375 0.625	1 1 1	0.375 0.375 0.625	321.5% 321.5% 520%
Cornelissen et al. (2018)	School readiness	Child care att.	Linear MTE	-0.1	0.22	0.08	0.082	0.975	813%
Kowalski (2023)	ER visits	HI coverage	Linear MTE	0.47	-0.58	0.2625	2.63	0.1	83.3%
Ito et al. (2023)	Electricity usage	Dynamic pricing	Local IV	-750	100	212.5	250	0.85	708.3%

Notes: Our (non-standardized) "minimal significant effect size" (MSE) with 95% (90%) confidence is 1.96 (1.64) times the standard error on the linearized interaction reported in Table 4, i.e., 0.203. Hence, the MSE yields 0.4 (0.33). Our dependent variable (the word recall score) has a standard deviation of 3.37, yielding a standardized MSE of 0.4/3.37 = 0.12. Note that this minimal significant effect size differs from the minimal detectable effect size for power calculations because we condition on the test result. The minimal detectable effect size is typically defined as the minimal effect detected as significant at the 5 % level in 80 % of all cases. This measure, however, cannot be computed without further assumptions. Our target parameter is the average difference in the effects of education on cognitive abilities between the fifth and the first quintile (i.e., four quintiles). The MTE maps the impact of a treatment for every quintile of the unobserved heterogeneity in the treatment choice. We approximate the average effect heterogeneity between quintiles of this unobserved heterogeneity by computing $\frac{MTE(p=0.1) - MTE(0.9)}{5-1}$ to compare the genetic heterogeneity to those of MTE studies (a perfect approximation would require identification at infinity and computing the integral between 0 and 0.2 and 0.8 and 1 for the first and fifth quintile, respectively). The reported values for $MTE(p = 0.1)$ and $MTE(p = 0.9)$ are approximated through eyeballing the corresponding MTE graphs. Nybohm (2017) interprets the degree of heterogeneity as low, justifying this finding with a low general heterogeneity in wages in Sweden. Further note that standard errors may be approximated, e.g., by averaging standard errors between treatment and control samples. The reported standard deviation of Ito et al. (2023) is approximated based on the average pre-intervention daily electricity consumption (the outcome measures hourly consumption during peak hours). Therefore, the reported SD may likely be a lower bound. See Møglstad and Torgovitsky (2024) for more applications of the MTE. Norm. sel. mod. = Normal selection model (maximum likelihood based functional form assumption (joint normality). HI coverage = Health insurance coverage. ER visits = Emergency room visits, College ed. = College education.

Table D.9: Robustness of the linearized 2SLS estimate

	DV: Total word recall score	
	Linearized $G \times E$ coefficient (1)	Standard error (2)
Baseline (main result, Table 4)	0.098	(0.174)
Birth cohort interacted with G_i	0.121	(0.425)
Fully interacted	0.044	(0.410)
Bandwidth 8 years	0.128	(0.191)
Bandwidth 5 years	0.121	(0.245)
Controls	Yes	
Observations	11,027	

Notes: This table presents robustness checks of the 2SLS estimation using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. We only present our main 2SLS result, the linearized $G \times E$ estimate, that represents a line through the $G \times E$ coefficients of the lowest and highest PGI quintile. This is done to compare average effects easily across methods. For reference, we report the effect from our main results. For “Birth cohort interacted with G_i ”, we add interactions between G_i and cohort trend t as well as $G_i \times t \times Z_i$ as controls on top of the baseline specification (Z_i being the instrument). For the fully interacted model, the controls include all baselines and interactions between t , G_i , Z_i , and X , where X is a vector of controls that includes gender, the principal components of the genetic data, parental education and survey wave fixed effects. Standard errors clustered at the individual level shown are in parentheses. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

Figure D.1: First Stage and reduced form by G_i

Notes: This figure shows a visualization of the first stage (upper panel) and the reduced form (lower panel) results of our regression discontinuity design by quintiles of the education PGI G_i using our main sample collapsed to cohort $\times G_i$ -averages. Dots correspond to G_i -specific sample means, lines are linear fits that are allowed to vary by Z_i . Regression results leading to this figure shown in Table D.10.

Table D.10: Estimates corresponding to Figure D.1

	Dependent variable:	
	Cohort average left school at 15 or later (E_i) (1)	Cohort average word recall (Y_i) (2)
$G_i = 1$	0.269 (0.032)***	8.933 (0.205)***
$G_i = 2$	0.415 (0.032)***	9.231 (0.205)***
$G_i = 3$	0.489 (0.032)***	9.493 (0.205)***
$G_i = 4$	0.503 (0.032)***	9.786 (0.205)***
$G_i = 5$	0.565 (0.032)***	9.866 (0.205)***
$Z_i \times (G_i = 1)$	0.667 (0.045)***	-0.284 (0.290)
$Z_i \times (G_i = 2)$	0.533 (0.045)***	0.054 (0.290)
$Z_i \times (G_i = 3)$	0.472 (0.045)***	-0.105 (0.290)
$Z_i \times (G_i = 4)$	0.430 (0.045)***	-0.222 (0.290)
$Z_i \times (G_i = 5)$	0.399 (0.045)***	0.244 (0.290)
t	0.007 (0.004)*	0.196 (0.024)***
$t \times Z_i$	-0.005 (0.005)	-0.039 (0.034)
Observations	100	100

Notes: This table presents estimates that correspond to averages and linear fits visualized in Figure D.1. This analysis uses our main sample collapsed to cohort $\times G_i$ -averages, yielding 100 observations (20 birth cohorts \times 5 gene groups). Columns (1) and (2) are OLS regression of the average proportion of the sample who left school at 15 or above (E_i) and average word recall score (Y_i), respectively, on G_i , interactions between G_i and the instrument Z_i , as well as a linear cohort trend t and the interaction of this trend with Z_i . We do not include a constant. Therefore, coefficients of G_i can be interpreted as sample means of the respective outcome variable in each quintile of G_i when $Z_i = 0$, i.e. on the left side of the cutoff. Interactions with Z are changes of these sample means for $Z_i = 1$, i.e. the right side. The coefficients of t and $t \times Z_i$ can be interpreted as the slope of the linear cohort trend for $Z_i = 0$ and its change for $Z_i = 1$. * $p < 0.1$, ** $p < 0.05$, and *** $p < 0.01$.

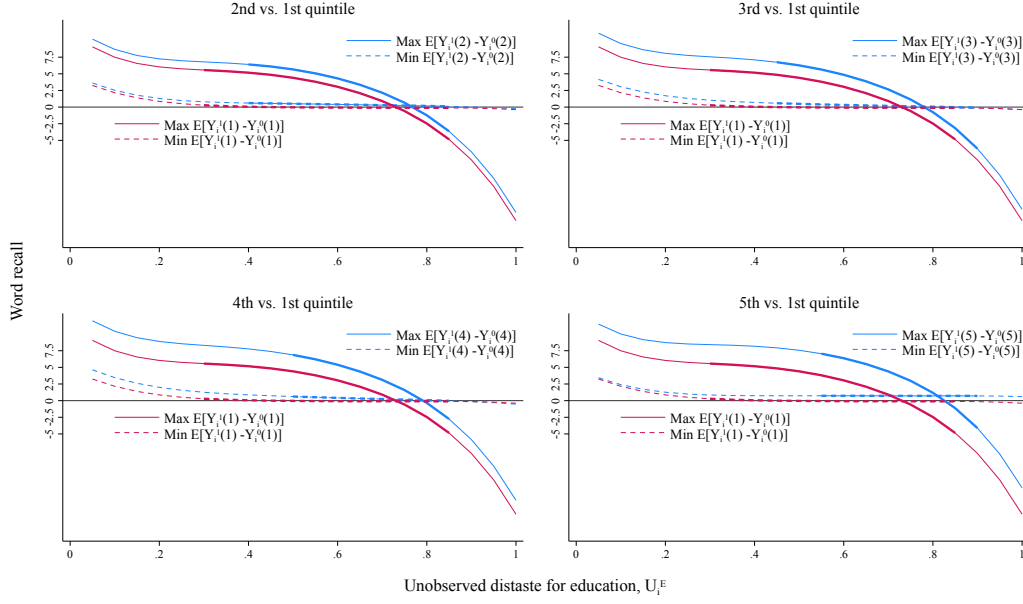


Figure D.2: Quintile comparisons of the interaction effect without never-takers

Notes: This figure shows the quintile comparisons of the interaction effect from Figure 7 when never-takers (their sample moments) are not used to construct the MTE bounds using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. For every PGI quintile, we estimate bounds: maxima (solid lines) and minima (dashed lines) at which the interaction effect is maximized/minimized. The bounds for quintiles 2–4 (in blue) are compared to those of the bottom quintile (in red), our reference category, yielding four comparisons. The gene-environment interaction is the difference between the blue and red curves at $U_i^E \in [0.55, 0.85]$. The thick part of the curves indicates the size of the complier share and its location on the U_i^E scale, both of which differ by PGI quintile.

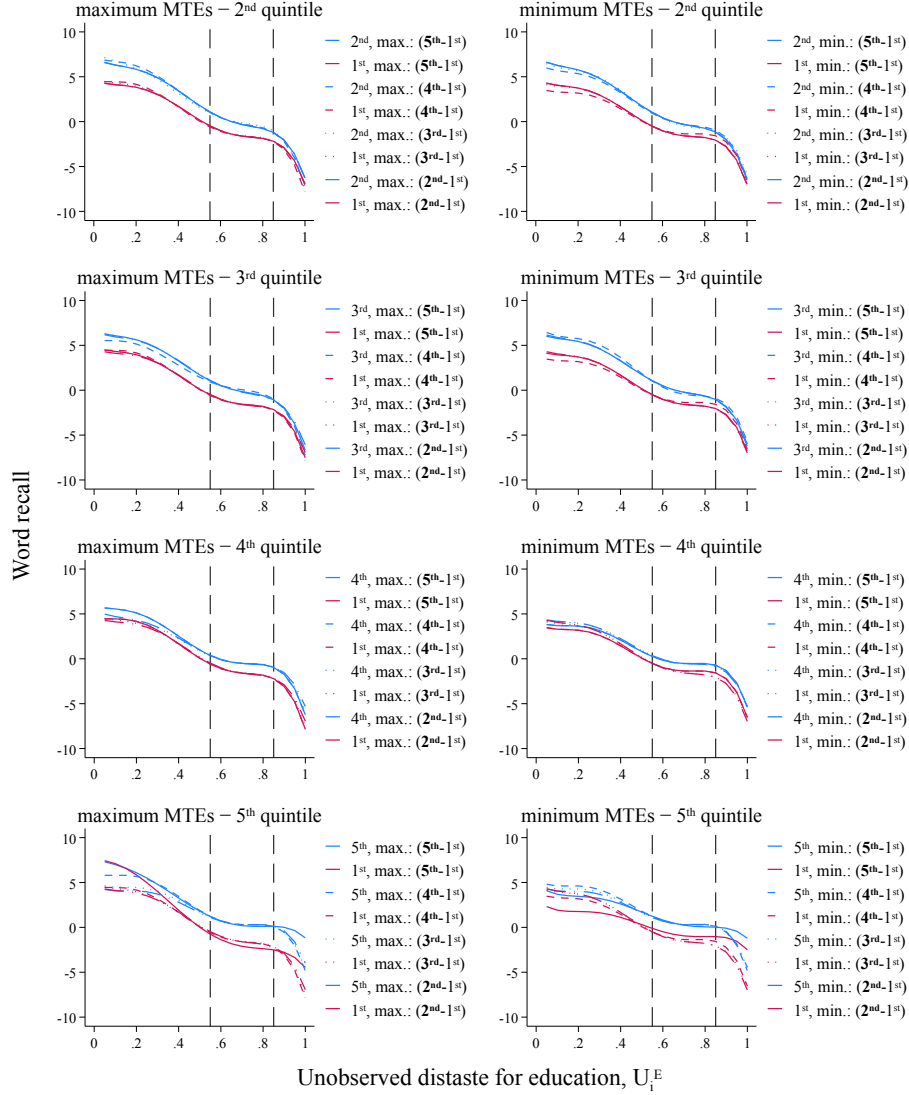


Figure D.3: MTEs when the target $G \times E$ parameter is adjusted to specific quintiles

Notes: This figure shows robustness checks for our main result in Figure 7 using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. Here, we optimize the interaction effect for different comparisons. Whereas our preferred specification optimizes the difference between the first and the fifth quintile (see Eq. 8), we generalize this approach and optimize differences between the first any other quintile such that $\beta_{G \times E}(0.55, 0.85, g) = \frac{1}{g-1} \int_{0.55}^{0.85} [m^1(u, g) - m^0(u, g)] - [m^1(u, 1) - m^0(u, 1)] du \quad \forall g \in \{2, 3, 4, 5\}$. The solid lines correspond to optimizing $g = 5$, our main result. The dashed lines show the optimization for $g = 4$, the dotted for $g = 3$, and the dashed-dotted line for $g = 2$. The respective quintile G_i used for the target parameter $\beta_{G \times E}(0.55, 0.85, g)$ is highlighted in bold. Maximized and minimized MTEs are shown separately, maximized MTEs in the left and minimized MTEs in the right column. The rows present pairwise comparisons between the first and another PGI quintile (the second quintile in the first row, the third in the second row, ...).

E What 2SLS is estimating

E.1 Formal derivation

Assume that E_i and G_i are binary. There are four potential outcomes $Y_i^j(G)$, $j \in \{0, 1\}$, $G \in \{0, 1\}$ of individual i . Only one is observed. The observation rule is

$$\begin{aligned}
 Y_i &= E_i \cdot G_i \cdot Y_i^1(1) + E_i \cdot (1 - G_i) \cdot Y_i^1(0) + (1 - E_i) \cdot G_i \cdot Y_i^0(1) + (1 - E_i) \cdot (1 - G_i) \cdot Y_i^0(0) \\
 &= Y_i^0(0) \\
 &\quad + (Y_i^1(0) - Y_i^0(0)) E_i \\
 &\quad + (Y_i^0(1) - Y_i^0(0)) G_i \\
 &\quad + (Y_i^1(1) - Y_i^0(1) - [Y_i^1(0) - Y_i^0(0)]) E_i \cdot G_i
 \end{aligned}$$

The second equation is the individual potential-outcome representation of the workhorse interaction model

$$Y_i = \beta_0 + \beta_1 E_i + \beta_2 G_i + \beta_3 G_i \times E_i + \varepsilon_i$$

Expressing this interaction equation as separate regressions for $G_i = 0$ and $G_i = 1$ yields

$$Y_i = \beta_0 + \beta_1 E_i + \varepsilon \quad \text{for } G_i = 0$$

$$Y_i = (\beta_0 + \beta_2) + (\beta_1 + \beta_3) E_i + \varepsilon \quad \text{for } G_i = 1$$

Environment E_i is often a choice variable, therefore endogenous and instrumented by Z_i , a binary instrument. In Wald notation, separately estimating 2SLS regressions for $G_i = 0$ and

$G_i = 1$ yields:

$$\begin{aligned}\widehat{\beta}_1 &= \frac{\mathbb{E}[Y_i|Z_i = 1, G_i = 0] - \mathbb{E}[Y_i|Z_i = 0, G_i = 0]}{\mathbb{E}[E_i|Z_i = 1, G_i = 0] - \mathbb{E}[E_i|Z_i = 0, G_i = 0]} && \text{for } G_i = 0 \\ \widehat{\beta}_1 + \widehat{\beta}_3 &= \frac{\mathbb{E}[Y_i|Z_i = 1, G_i = 1] - \mathbb{E}[Y_i|Z_i = 0, G_i = 1]}{\mathbb{E}[E_i|Z_i = 1, G_i = 1] - \mathbb{E}[E_i|Z_i = 0, G_i = 1]} && \text{for } G_i = 1\end{aligned}$$

Using the LATE theorem ([Imbens and Angrist, 1994](#)—2SLS estimates are average treatment effects for the compliers), we can rewrite these expressions as:

$$\begin{aligned}\widehat{\beta}_1 &= \mathbb{E}[Y_i^1(0) - Y_i^0(0)|C(G_i = 0)] \\ \widehat{\beta}_1 + \widehat{\beta}_3 &= \mathbb{E}[Y_i^1(1) - Y_i^0(1)|C(G_i = 1)]\end{aligned}$$

The mechanics of the LATE require that the group-specific effects ($\widehat{\beta}_1$ and $\widehat{\beta}_1 + \widehat{\beta}_3$) are average treatment effects for the G_i -specific compliers. Without further covariates, the joint interaction regression specification is as flexible as the separate ones. The mechanics of interaction models attribute any difference in the causal effects of E_i on Y_i between $G_i = 0$ and $G_i = 1$ to the interaction coefficient. In essence, the interaction model is numerically identical to separate estimations. In the interaction model, any difference between the G_i -specific LATEs is mechanically attributed to $\widehat{\beta}_3$. Hence, using the expressions above, this difference amounts to:

$$\widehat{\beta}_3 = (\widehat{\beta}_1 + \widehat{\beta}_3) - \widehat{\beta}_1 = \mathbb{E}[Y_i^1(1) - Y_i^0(1)|C(G_i = 1)] - \mathbb{E}[Y_i^1(0) - Y_i^0(0)|C(G_i = 0)]$$

This demonstrates that the interaction coefficient reflects differences in G_i -specific LATEs.

E.2 Simulation model

To visualize that 2SLS is unable to disentangle interaction effects from shifts in complier groups, as discussed in Section 4, we set up a simple simulation model. Assume the following arbitrary parametrizations of the potential outcomes, where, for simplicity, we leave out observable variables X_i :

$$\begin{aligned} Y_i^1(1) &= 2.3 + \varepsilon_i^1, & Y_i^1(0) &= 0.5 + \varepsilon_i^1, & Y_i^0(1) &= 0.3 + \varepsilon_i^0, & Y_i^0(0) &= 0 + \varepsilon_i^0 \\ E_i &= \mathbb{1}\{0.23 + 2.5G_i - 4Z_i + 3Z_i \times G_i > -(\varepsilon_i^1 - \varepsilon_i^0)\} \\ Z_i, G_i &= \text{Bernoulli distributed with } p = 0.5, \end{aligned}$$

where $\varepsilon_i^0 = \varepsilon_i^0(G_i)$ and $\varepsilon_i^1 = \varepsilon_i^1(G_i)$ are error terms. Here, we make the simplifying assumption that $\varepsilon^1(1) = \varepsilon_i^1(0) = \varepsilon_i^1$ and $\varepsilon_i^0(1) = \varepsilon_i^0(0) = \varepsilon_i^0$. In this simulation — but not in the application later — they are assumed to follow a bivariate normal distribution with the following parameters:

$$\begin{pmatrix} \varepsilon_i^1 \\ \varepsilon_i^0 \end{pmatrix} \sim N \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 & 0.4 \\ 0.4 & 2 \end{pmatrix} \right].$$

Setting $\varepsilon^1(1) = \varepsilon_i^1(0) = \varepsilon_i^1$ and $\varepsilon_i^0(1) = \varepsilon_i^0(0) = \varepsilon_i^0$ does not affect the main line of argumentation and is merely for a simple exposition. It restricts the gene-environment effect to 1.5 for each individual and, thus, each complier type. Assuming four different error terms allows for a different gene-environment interaction effect by complier type. Our argument is not affected by that, and neither does our solution need this restriction, nor do we make this assumption in the application in Sections 2 to 5.

E.3 Estimating outcome means for always-takers, compliers, and never-takers following Imbens and Rubin (1997)

Imbens and Rubin (1997) show that with a binary treatment indicator E and a binary instrument Z , four potential outcomes can be identified for always-takers (AT), compliers (C), and never-takers (NT): $\mathbb{E}[Y_i^1|AT]$, $\mathbb{E}[Y_i^1|C]$, $\mathbb{E}[Y_i^0|C]$, and $\mathbb{E}[Y_i^0|NT]$.

To estimate the four numbers, group individuals into four different strata, based on their observed E_i and Z_i .

$$\mathbb{E}[Y_i|E_i = 1, Z_i = 0] = \mathbb{E}[Y_i^1|AT, Z_i = 0] = \mathbb{E}[Y_i^1|AT]$$

$$\mathbb{E}[Y_i|E_i = 1, Z_i = 1]$$

$$\mathbb{E}[Y_i|E_i = 0, Z_i = 1] = \mathbb{E}[Y_i^0|NT, Z_i = 1] = \mathbb{E}[Y_i^0|NT]$$

$$\mathbb{E}[Y_i|E_i = 0, Z_i = 0]$$

Under the assumption that there are no defiers (monotonicity), individuals with $E_i = 1, Z_i = 0$ must be always-takers, as they take the treatment even when the instrument is switched off ($Z_i = 0$). Conversely, individuals with $E_i = 0, Z_i = 1$ reveal themselves as never-takers, since they do not take the treatment even when the instrument is switched on. Therefore, $\mathbb{E}[Y_i^1|AT, Z_i = 0]$ and $\mathbb{E}[Y_i^0|NT, Z_i = 1]$ can directly be estimated from the data. If the instrument is exogenous, this also provides consistent estimates of $\mathbb{E}[Y_i^1|AT]$ and $\mathbb{E}[Y_i^0|NT]$.

The group with $E_i = 1, Z_i = 1$ is a mixture of compliers and always-takers. The group with $E_i = 0, Z_i = 0$ is a mixture of compliers and never-takers. Using the law of total

probability we can express the estimable number $\mathbb{E}[Y_i|E_i = 1, Z_i = 1]$ as follows

$$\begin{aligned} \mathbb{E}[Y_i|E_i = 1, Z_i = 1] &= \mathbb{E}[Y_i|E_i = 1, Z_i = 1, AT] \cdot Pr(AT|E_i = 1, Z_i = 1) \\ &\quad + \mathbb{E}[Y_i|E_i = 1, Z_i = 1, C] \cdot Pr(C|E_i = 1, Z_i = 1) \end{aligned} \quad (\text{E.5})$$

where $Pr(AT|E_i = 1, Z_i = 1)$ and $Pr(C|E_i = 1, Z_i = 1)$ are the shares of AT and C within this stratum.

Further, it holds that

- $\mathbb{E}[Y_i|E_i = 1, Z_i = 1, AT] = \mathbb{E}[Y_i|E_i = 1, Z_i = 0, AT] = \mathbb{E}[Y_i|E_i = 1, Z_i = 0]$ using instrument exogeneity and the reasoning for AT used two paragraphs above.
- $\mathbb{E}[Y_i|E_i = 1, Z_i = 1, C] = \mathbb{E}[Y_i^1|Z_i = 1, C] = \mathbb{E}[Y_i^1|C]$
- $Pr(AT|E_i = 1, Z_i = 1) = \frac{\phi_{AT}}{\phi_{AT} + \phi_C}$ we ϕ_{AT} denotes the share of always-takers in the sample and ϕ_C the share of compliers
- $Pr(C|E_i = 1, Z_i = 1) = \frac{\phi_C}{\phi_{AT} + \phi_C}$

Now, plugging in these four expressions into Equation (E.5) and solving for $\mathbb{E}[Y_i^1|C]$ yields:

$$\mathbb{E}[Y_i^1|C] = \mathbb{E}[Y_i|E_i = 1, Z_i = 1] \frac{\phi_{AT} + \phi_C}{\phi_C} - \mathbb{E}[Y_i|E_i = 1, Z_i = 0] \frac{\phi_{AT}}{\phi_C} \quad (\text{E.6})$$

The average treated outcome for compliers is now expressed entirely in terms of the estimable means and type shares defined above, allowing for a straightforward calculation. Similar

reasoning gives:

$$\begin{aligned}
 \mathbb{E}[Y_i|E_i = 0, Z_i = 0] &= \mathbb{E}[Y_i|E_i = 0, Z_i = 0, NT] \cdot Pr(NT|E_i = 0, Z_i = 0) \\
 &\quad + \mathbb{E}[Y_i|E_i = 0, Z_i = 0, C] \cdot Pr(C|E_i = 0, Z_i = 0) \\
 &= \mathbb{E}[Y_i|E_i = 0, Z_i = 1] \cdot \frac{\phi_{NT}}{\phi_{NT} + \phi_C} + \mathbb{E}[Y_i^0|C] \cdot \frac{\phi_C}{\phi_{NT} + \phi_C}
 \end{aligned}$$

Solving for $\mathbb{E}[Y_i^0|C]$ gives:

$$\mathbb{E}[Y_i^0|C] = \mathbb{E}[Y_i|E_i = 0, Z_i = 0] \frac{\phi_{NT} + \phi_C}{\phi_C} - \mathbb{E}[Y_i|E_i = 0, Z_i = 1] \frac{\phi_{NT}}{\phi_C} \quad (\text{E.7})$$

We now discuss how the group shares ϕ_{AT} , ϕ_{NT} , ϕ_C can be estimated. First, we calculate $Pr(E_i = 1|Z_i = 0)$, which is equivalent to the sample average of E_i among individuals with $Z_i = 0$. If the instrument is random, this is equal to the share of always-takers with $Z = 1$ and the share of always-takers in the full population. Analogously, we can calculate $Pr(E_i = 0|Z_i = 1)$ as the sample average of $1 - E_i$ among individuals with $Z_i = 1$. This provides the share of never-takers with $Z_i = 1$, which is equal to the share of never-takers with $Z = 0$ and in the full population. Since there are only three types and their shares must add up to one, we can calculate the share of compliers simply as $\phi_C = 1 - \phi_{AT} - \phi_{NT}$. Alternatively, we can obtain these shares by running regression $E_i = \beta_0 + \beta_1 Z_i + \varepsilon$. The share of always-takers is given by the intercept β_0 and the share of compliers by the regression coefficient of Z , β_1 . The share of never-takers can then be calculated as $\phi_{NT} = 1 - \beta_0 - \beta_1$.

Now, we have expressed all quantities we can obtain using means and shares directly observable in the data. The most straightforward way to obtain all observable means required

is to run a single regression of the outcome Y on indicator variables for each subgroup:

$$\begin{aligned}
 Y = & \gamma_{0,0} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 0\} \\
 & + \gamma_{1,0} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 0\} \\
 & + \gamma_{0,1} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 1\} \\
 & + \gamma_{1,1} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 1\} + \varepsilon
 \end{aligned} \tag{E.8}$$

The coefficients from this regression correspond directly to the sample means as follows:

$$\begin{aligned}
 \mathbb{E}[Y_i | E_i = 1, Z_i = 0] &= \gamma_{1,0} && \text{AT} \\
 \mathbb{E}[Y_i | E_i = 1, Z_i = 1] &= \gamma_{1,1} && \text{AT+C} \\
 \mathbb{E}[Y_i | E_i = 0, Z_i = 1] &= \gamma_{0,1} && \text{NT} \\
 \mathbb{E}[Y_i | E_i = 0, Z_i = 0] &= \gamma_{0,0} && \text{NT+C}
 \end{aligned}$$

Since in our setting we want to obtain not four outcome means, but 20 — four for each of the five quintiles of G_i — we interact Eq. E.8 with G_i and run the following regression:

$$\begin{aligned}
Y_i = & \delta_{1,0,0} \mathbb{1}\{G_i = 1\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 0\} + \delta_{1,1,0} \mathbb{1}\{G_i = 1\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 0\} \\
& + \delta_{1,0,1} \mathbb{1}\{G_i = 1\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 1\} + \delta_{1,1,1} \mathbb{1}\{G_i = 1\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 1\} \\
& + \delta_{2,0,0} \mathbb{1}\{G_i = 2\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 0\} + \delta_{2,1,0} \mathbb{1}\{G_i = 2\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 0\} \\
& + \delta_{2,0,1} \mathbb{1}\{G_i = 2\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 1\} + \delta_{2,1,1} \mathbb{1}\{G_i = 2\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 1\} \\
& + \delta_{3,0,0} \mathbb{1}\{G_i = 3\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 0\} + \delta_{3,1,0} \mathbb{1}\{G_i = 3\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 0\} \\
& + \delta_{3,0,1} \mathbb{1}\{G_i = 3\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 1\} + \delta_{3,1,1} \mathbb{1}\{G_i = 3\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 1\} \quad (\text{E.9}) \\
& + \delta_{4,0,0} \mathbb{1}\{G_i = 4\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 0\} + \delta_{4,1,0} \mathbb{1}\{G_i = 4\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 0\} \\
& + \delta_{4,0,1} \mathbb{1}\{G_i = 4\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 1\} + \delta_{4,1,1} \mathbb{1}\{G_i = 4\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 1\} \\
& + \delta_{5,0,0} \mathbb{1}\{G_i = 5\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 0\} + \delta_{5,1,0} \mathbb{1}\{G_i = 5\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 0\} \\
& + \delta_{5,0,1} \mathbb{1}\{G_i = 5\} \mathbb{1}\{E_i = 0\} \mathbb{1}\{Z_i = 1\} + \delta_{5,1,1} \mathbb{1}\{G_i = 5\} \mathbb{1}\{E_i = 1\} \mathbb{1}\{Z_i = 1\} \\
& + \text{controls} + \zeta_i
\end{aligned}$$

To obtain the shares of always-takers, compliers and never-takers for each quintile of G_i , we run the following second regression:

$$\begin{aligned}
E_i = & \pi_{1,0} \mathbb{1}\{G_i = 1\} \mathbb{1}\{Z_i = 0\} + \pi_{1,1} \mathbb{1}\{G_i = 1\} \mathbb{1}\{Z_i = 1\} \\
& + \pi_{2,0} \mathbb{1}\{G_i = 2\} \mathbb{1}\{Z_i = 0\} + \pi_{2,1} \mathbb{1}\{G_i = 2\} \mathbb{1}\{Z_i = 1\} \\
& + \pi_{3,0} \mathbb{1}\{G_i = 3\} \mathbb{1}\{Z_i = 0\} + \pi_{3,1} \mathbb{1}\{G_i = 3\} \mathbb{1}\{Z_i = 1\} \\
& + \pi_{4,0} \mathbb{1}\{G_i = 4\} \mathbb{1}\{Z_i = 0\} + \pi_{4,1} \mathbb{1}\{G_i = 4\} \mathbb{1}\{Z_i = 1\} \\
& + \pi_{5,0} \mathbb{1}\{G_i = 5\} \mathbb{1}\{Z_i = 0\} + \pi_{5,1} \mathbb{1}\{G_i = 5\} \mathbb{1}\{Z_i = 1\} \\
& + \text{controls} + \tau_i
\end{aligned} \tag{E.10}$$

Based on these estimates, we apply the Equations (E.6) and (E.7) to compute outcome means for treated and untreated compliers for each G_i :

$$\begin{aligned}
\mathbb{E}[Y_i^1(G)|C, G_i = g] &= \frac{\delta_{g,1,1}\pi_{g,1} - \delta_{g,1,0}\pi_{g,0}}{\pi_{g,1} - \pi_{g,0}} \\
\mathbb{E}[Y_i^0(G)|C, G_i = g] &= \frac{\delta_{g,0,0}\pi_{g,0} - \delta_{g,0,1}\pi_{g,1}}{\pi_{g,1} - \pi_{g,0}} \\
\mathbb{E}[Y_i^0(G)|NT, G_i = g] &= \delta_{g,0,1} \\
\mathbb{E}[Y_i^1(G)|AT, G_i = g] &= \delta_{g,1,0}
\end{aligned}$$

Treated means for always-takers and untreated means for never-takers are obtained directly.

These 20 means are plotted in Figure F.1.

F Details on the MTE estimation

F.1 A brief introduction to MTEs

We start by briefly summarizing the classic MTE framework by [Heckman and Vytlacil \(2005\)](#), adjusted to our notation. See [Heckman and Vytlacil \(2005\)](#) for an extensive introduction to MTEs, their derivation, and traditional ways to estimate them with continuous instruments.

Assume that the potential outcomes of individual i are defined by the following functions: $Y_i^j(G_i) = \mu^j(G_i, X_i) + \varepsilon_i^j(G_i)$, $j \in \{0, 1\}$, $G_i \in \{0, 1\}$, where j denotes potential outcomes of educational status and G_i the attributes of genetic endowment. $\mu^j(G_i, X_i)$ is a function of genetic endowment G_i and observable characteristics X_i , and $\varepsilon_i^j(G_i)$ is an unobservable part.

We model the choice E_i in a generalized Roy framework ([Roy, 1951](#)), where individuals choose E_i if the (expected) returns to education exceed monetary and/or non-monetary costs $C_i = \mu^C(G_i, X_i, Z_i) + U_i^C$. Costs depend on G_i , the observable characteristics X_i , an instrumental variable Z_i and an unobservable term U_i^C . Note that Z_i does not directly affect $Y_i^j(G_i)$. The decision rule for E_i (depending on the realization of $G_i = g$) reads:

$$\begin{aligned} E_i(G_i) = 1 &\Leftrightarrow Y_i^1(G_i) - Y_i^0(G_i) > C_i \\ &\Leftrightarrow \mu^1(G_i, X_i) - \mu^0(G_i, X_i) - \mu^C(G_i, X_i, Z_i) > -(\varepsilon_i^1(G_i) - \varepsilon_i^0(G_i) - U_i^C(G_i)) \\ &\Leftrightarrow \mu^E(G_i, X_i, Z_i) > V_i(G_i) \end{aligned}$$

While not necessary for any theoretical result, $\mu^E(G_i, X_i, Z_i) = \mu^1(G_i, X_i) - \mu^0(G_i, X_i) - \mu^C(G_i, X_i, Z_i)$ can be represented as a linear index, such as:

$$\mu^E(G_i, X_i, Z_i) = \pi_0 + \pi_1 G_i + \pi_2 Z_i + \pi_3 Z_i \cdot G_i + \pi X_i + V_i(G_i)$$

where $V_i(G_i) = -(\varepsilon_i^1(G_i) - \varepsilon_i^0(G_i) - U_i^C)$ is the unobservable term. The decision rule implies that E_i correlates with $\varepsilon_i^1(G_i)$ and $\varepsilon_i^0(G_i)$ and, thus, $V_i(G_i)$, which renders E_i endogenous.

In the spirit of [Heckman and Vytlačil \(2005\)](#) we rewrite the choice equation as:

$$\begin{aligned} E(G_i) &= \mathbb{1}\{\mu^E(G_i, X_i, Z_i) \geq V_i(G_i)\} \\ &= \mathbb{1}\{F_V(\mu^E(G_i, X_i, Z_i)) \geq F_V(V_i(G_i))\} \\ &= \mathbb{1}\{\Pr(V_i(G_i) \leq \mu^E(G_i, X_i, Z_i)) \geq F_V(V_i(G_i))\} \\ &= \mathbb{1}\{\Pr(E(G_i = 1)|X_i, Z_i) \geq U_i^E\} \\ &= \mathbb{1}\{PS(G_i, X_i, Z_i) \geq U_i^E\} \end{aligned}$$

The second step applies a monotonic transformation $F_V(\cdot)$ — which is the cumulative density of $V_i(G)$ — to both sides of the inequality. $F_V(\cdot)$ evaluated at the point $\mu^E(G_i, X_i, Z_i)$ is defined as $\Pr(V_i(G) \leq \mu^E(G_i, X_i, Z_i))$ and, referring to the choice equation, the same as $\Pr(E(G_i = 1)|X_i, Z_i)$. This choice probability based on observable characteristics is the propensity score, and we abbreviate it by $PS(G_i, X_i, Z_i)$. Irrespective of the underlying distribution of $V_i(G)$, the unobserved term U_i^E is uniformly distributed on the unit interval and comprises the unobserved heterogeneity correlating with the decision to take E_i . Low values of unobserved resistance to more education U_i^E increase $PS(G_i, X_i, Z_i)$, leading to

$E_i = 1$. This corresponds to high unobserved preferences for E_i , whereas large values of U_i^E indicate a high distaste for E_i .

MTEs are estimates of the causal effect of education on the outcome Y_i at certain values of $U_i^E = u$. That is, $\mathbb{E}[Y^1(G) - Y^0(G)|U_i^E = u]$. The MTEs are identified by those individuals who, at $U_i^E = u$, are indifferent between choosing $E_i = 0$ and $E_i = 1$. Referring to the choice equation, this is the group for whom the realization p of the propensity score $PS(G_i, X_i, Z_i) = p = u$. For our framework, the quantities $\mathbb{E}[Y_i^1(G)|U_i^E = u]$ and $\mathbb{E}[Y_i^0(G)|U_i^E = u]$ are essential (as their difference is the MTE). We follow the literature and call these quantities marginal treatment response curves (MTRs).

Note that conditioning on a narrow range of the propensity score (PS) and then estimating Eq. (6) in this more homogeneous subsample only works when both PS and Z are continuous.²⁰ In our setting, the propensity score is discrete. Evaluating the compulsory-schooling cutoff that defines the instrument on either side for each genetic stratum (G) yields variation in the PS that identifies G -specific always-takers, compliers, and never-takers. Hence, within G , the binary instrument generates only two PS values, providing insufficient variation to make propensity scores comparable across different values of G . For this reason, we estimate marginal treatment effects instead. MTEs recover how treatment effects vary with unobserved resistance to treatment. Under identifying restrictions (e.g., our constraints), the MTE approach does not require full propensity-score overlap across covariates (G).

²⁰In this case, one can partial out observables and apply a local linear regression—the standard semiparametric approach to identifying the MTE. With a discrete instrument, however, conditioning on PS and estimating Eq. (6) by 2SLS will typically drop observations based on realizations of Z , making the conditional sample endogenous.

F.2 Options for MTE estimation

There are many different ways to estimate MTRs and MTEs, depending on the underlying data, setting (e.g., continuous or binary instrumental variables) and the assumptions the researcher wants to impose (e.g., functional form assumptions for the MTE, separability between observed and unobserved terms). In our case, with a binary instrument, there are at least three options.

1. Linear MTEs: Estimate different expected values of potential outcomes $\mathbb{E}[Y_i^1(G)|AT]$, $\mathbb{E}[Y_i^1(G)|C]$, $\mathbb{E}[Y_i^0(G)|C]$, and $\mathbb{E}[Y_i^0(G)|NT]$ (see [Imbens and Rubin, 1997](#)) for each value of G_i . Plotted on the U_i^E unit interval and assuming linearity, we can fit lines through each pair of points, one for treated and one for untreated potential outcomes ([Brinch et al., 2017](#)), which provide the MTRs. The linearity ensures that the lines run through the respective (type-specific) midpoints on the U_i^E scale. The difference between the two MTR lines is the linear MTE by G_i , and the four differences between the five G_i -specific MTEs inform about the interaction effects.
2. Semi-parametric estimation: Relax the linearity assumption but impose additive separability between controls X_i and error terms. That is, specify the potential outcomes $Y_i^j(G_i) = \mu^j(G_i, X_i) + \varepsilon_i^j(G_i)$, as we have already done above, instead of the more general form $Y_i^j(G_i) = f^j(G_i, X_i, \varepsilon_i^j)$ with some arbitrary function f . This allows variation in X_i to parametrically or semi-parametrically identify the MTEs, since a binary instrument alone cannot provide this ([Brinch et al., 2017](#)).
3. Bounds around flexible polynomial shapes: Allow for a wide range of flexible polynomial shapes of the MTEs and subsequently restrict the shapes. This can be achieved

by requiring the curves to reproduce observable sample analogs and imposing further reasonable assumptions derived from theory and the data. The target parameter the researcher aims to identify can be bounded by the two shapes that produce minimum and maximum values (Mogstad *et al.*, 2018).

A linearity assumption is hard to justify a priori. Furthermore, although additive separability is commonly assumed across the entire literature that uses regression models, we do not benefit from it for a semi-parametric identification of the MTEs. This is because we only use a sparse set of control variables that do not add sufficient variation in the propensity score, which would help identify substantially more than the four points from the first approach. Overall, the third approach appears to be the most suitable for our setting.

F.3 Linear MTEs as a benchmark

As a benchmark, we estimate linear MTEs according to Brinch *et al.* (2017). They help to illustrate the setting and show general trends. They are also informative about underlying shape restrictions.

To achieve estimates of linear MTEs, we begin by estimating type-specific expected outcomes, that correspond to group means in the data: $\mathbb{E}[Y_i^1(G)|AT]$, $\mathbb{E}[Y_i^1(G)|C]$, $\mathbb{E}[Y_i^0(G)|C]$, $\mathbb{E}[Y_i^0(G)|NT]$. We also estimate the shares of AT, C, and NT for each quintile of the polygenic index. Appendix E.3 presents the details on generating these 35 values by applying the Imbens and Rubin (1997) method. We visualize the 20 means (circles and diamonds, depending on treatment status) as well as the 15 type shares (horizontal lines at the bottom) in Figure F.1.²¹ Again, we sort the three complier types on the unit interval, according to

²¹Within each of the five quintiles of G , we estimate four means: the treated mean for always-takers, the treated mean for compliers, the untreated mean for compliers, and the untreated mean for never-takers. This

their willingness to take education. Always-takers have the highest willingness and are located at the left. For example, the share of always-takers in the lowest PGI quintile (lowest horizontal line) is 22 %. The share of compliers with $G_i = 1$ is 68 %. They are located between 0.22 and 0.9 on the U_i^E unit interval. The remaining 10 % are never-takers. Following [Kowalski \(2023\)](#), we use the midpoints of the range where each type is located to place the potential outcomes (circles and diamonds) on the x-axis. At the same time, the y-axis measures the size of the estimated potential outcomes. The blue circles denote treated potential outcomes $\mathbb{E}[Y_i^1(G)]$ while the red diamonds denote untreated potential outcomes $\mathbb{E}[Y_i^0(G)]$. The numbers next to the markers refer to the realization of G_i .

The lines through the points produce MTRs under the linearity assumption, allowing us to identify them using the two points. In principle, the lines can be extrapolated to the full unit interval and taking differences between $\mathbb{E}[Y_i^1(G)|U_i^E = u]$ and $\mathbb{E}[Y_i^0(G)|U_i^E = u]$ would yield the MTEs by G_i . Comparing the resulting five MTEs provides insight into gene-environment interactions. However, as mentioned above, the linearity assumption, which will drive the final results, is hard to defend a priori. Nevertheless, this analysis has important implications for our bounding approach to compute our main results (Section [5.1](#)). Treated potential outcomes (blue) are higher for always-takers than compliers, causing the $\mathbb{E}[Y_i^1(G)|U_i^E = u]$ -MTR curves to have a negative slope. Therefore, there seems to be a correlation between types and our dependent variable. Moreover, the $\mathbb{E}[Y_i^1(G)|U_i^E = u]$ -MTR curves are fairly parallel, with no substantial slope differences. They represent level shifts in treated outcomes by G_i . These results are used to justify constraints 3 and 5 in Section [5.1](#).

gives $4 \times 5 = 20$ means. In addition, we report the shares of always-takers, compliers, and never-takers in each quintile of G , yielding $3 \times 5 = 15$ type shares.

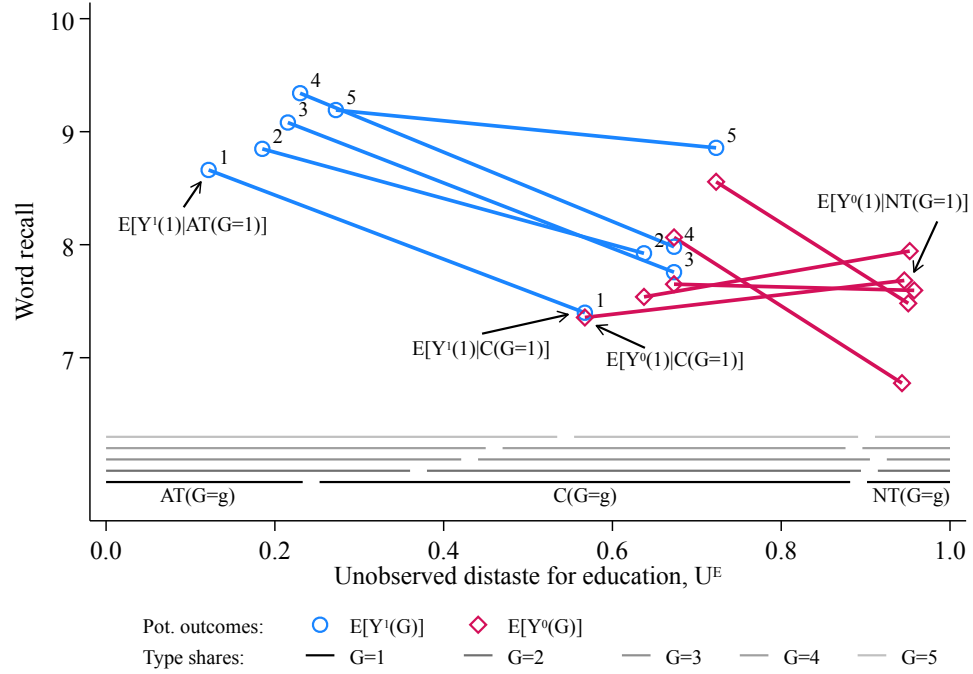


Figure F.1: Linear potential outcome curves

Notes: This figure shows the 20 estimated potential outcomes $E[Y_i^j(G)]$ using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. The lines through them represent linear MTRs. Red diamonds refers to potential outcomes for $E_i = 0$; blue circles to $E_i = 1$. Thus, for example, the red line labelled “1” shows our estimate of the potential outcome curve of $Y_i^0(1)$; the blue line labelled “3” shows the curve for $Y_i^1(3)$. Horizontal lines at the bottom show type shares by quintiles of the educational attainment PGI, with their locations on the unit interval in ascending order, starting with $G_i = 1$ (the lowest, black line) and $G_i = 5$ (the highest, lightest line). We provide detailed descriptions of potential outcomes as text with arrows for $G_i = 1$ (as an example and to maintain readability).

The picture is less clear for the untreated outcomes (red diamonds). Here, we see that the outcomes of the untreated compliers are, on average, slightly smaller than those of the treated compliers, suggesting positive effects of education on word recall. We can replicate the 2SLS finding of a zero effect for $C(G_i = 1)$ and positive effects for $C(G_i = 2)$ and $C(G_i = 5)$. However, the estimates for never-takers are less clear, as they are above those for the untreated compliers for the first and second quintiles, resulting in positive slopes of the two lowest $E[Y_i^0(G)|U_i^E = u]$ -MTRs. Given this ambiguous result and the small share of

never-takers in the data, we include a robustness check of our main result where we estimate MTEs without relying on never-takers in Section 5.4. However, connecting to the argument in [Clark and Royer \(2013\)](#), never-takers are the youngest in their class and, when leaving school. Hence, as it contains valid information, we opted to use never-takers in our main specification.

F.4 Details on the method by Mogstad et al. (2018)

Our approach assumes a parametric shape of the MTRs, which, however, is extremely flexible by specifying them as Bernstein polynomials.²² The Bernstein polynomials are defined as

$$\mathbb{E}[Y_i^j(G_i = g)|U_i^E = u, G_i = g] = \sum_{v=0}^n \theta_v^{jg} \binom{n}{v} u^v (1-u)^{n-v},$$

where u is a specific point on the unit interval, j refers to the treatment state, g is the PGI quintile, and n is the polynomial degree. We choose $n = 5$. Therefore, we have $n + 1 = 6$ parameters $(\theta_0^{jg}, \dots, \theta_n^{jg})$ that determine each of the 10 MTR functions (5 for $j = 0$ and 5 for $j = 1$). See Figure F.2 for a graphical representation of the Bernstein base functions.

The θ -parameters of the MTRs are derived from a linear programming exercise which can be described as follows: (i) Among all theoretically possible MTRs, consider only those that fulfill certain constraints. These constraints—combinations of assumptions and information from the data—are set by the researcher, and we lay them out in the main text. (ii) Among all MTRs that fulfill the constraints, find those that maximize the target parameter. (iii)

²²Again, we keep the explanation of the approach brief and largely verbal and refer to [Mogstad et al. \(2018\)](#) for a comprehensive treatment and technical exposition.

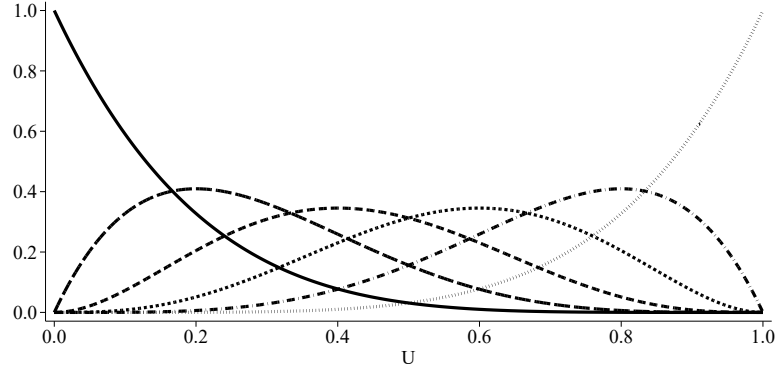


Figure F.2: Graphical representation of the Bernstein base functions

Notes: This figure depicts the six Bernstein base functions that compose a Bernstein polynomial of degree five in simulated data. The formula for each base function reads $b_{v,n}(u) = \binom{n}{v} u^v (1-u)^{n-v}$, where $n = 5$ is the degree, v denotes the specific base function and u is a specific grid point on the unit interval.

Among all MTRs that fulfill the constraints, find those that minimize the target parameter.

As a result, we get a bound around the target parameter.

In total, there are 60 parameters: 6 (determined by the polynomial degree) times 2 (treated and untreated cases) times 5 (different values of G_i). Estimating the bounds (i.e., choosing the 60 parameters) involves solving a linear programming problem where constraints on the Bernstein polynomial shapes can be represented as constraints on the parameters θ (Rose and Shem-Tov, 2021).

F.5 Estimation/computation protocol

Part A: Moment estimation by OLS

1. Estimate the complier shares using Eq. (E.10). Assemble the resulting $\pi_{g,z}$ into a vector, ordered by G and (within G) by E . Call this vector π .

2. Estimate outcome means for the complier types (i.e., the “moments”) using Eq. (E.9). For $Z = 1$ and $Z = 0$, assemble the resulting $\delta_{g,j,z}$ into a vector, ordered by G , (within G) by E , and (within G and Z) by Z . Call this vector δ .

Part B: Linear Programming Approach

1. Set up the basic ingredients (grid vector and Bernstein matrix)
 - (i) Choose grid points ($m = 20$),
 - (ii) Choose degree of Bernstein polynomials ($n = 5$),
 - (iii) Compute grid vector $\mathbf{u}' = \begin{pmatrix} \frac{1}{m} & \frac{2}{m} & \dots & 1 \end{pmatrix} \in \mathbb{R}^{1 \times m}$.
 - (iv) Compute Bernstein base vector as $\mathbf{b}_v = \binom{n}{v} \mathbf{u}^v (1 - \mathbf{u})^{n-v} \quad \forall v \in \{0, \dots, n\}$. Bind these base vectors to a matrix $\mathbf{b} = \begin{pmatrix} \mathbf{b}_0 & \dots & \mathbf{b}_n \end{pmatrix} \in \mathbb{R}^{(m \times n)}$.

Remark: The choices of m and n are the same as in [Rose and Shem-Tov, 2021](#)). Increasing the number of grid points does not significantly affect our results but substantially increases computation time.
2. Tailor this to our specific problem by expanding \mathbf{b} by $E = j$ states and $G = g$ strata:
 - (i) Use the Kronecker product to expand the polynomials for the treatment state E : $\mathbf{b}^E = \mathbf{I}(2) \otimes \mathbf{b}$. This transforms the $(m \times n)$ matrix \mathbf{b} into a $(2m \times 2n)$ matrix \mathbf{b}^E .

(ii) Use the Kronecker product to expand the \mathbf{b}^E matrix for the genetic strata

G : $\mathbf{b}^{G \times E} = \mathbf{I}(\mathbf{5}) \otimes \mathbf{b}^E$. This transforms the $(2m \times 2n)$ matrix \mathbf{b}^E into an $(10m \times 10n)$ matrix $\mathbf{b}^{E \times G}$.

3. Specify the weighting matrix for the target parameters, ω .

(i) For each element u_i of \mathbf{u} and treatment state j compute

$$\omega^*(j, u_i) = \begin{cases} -1^{(j+1)} & \text{if } u_i \in [0.55, 0.85] \\ 0 & \text{else.} \end{cases}$$

(ii) For treatment state j , bind all m elements to a $m \times 1$ vector

$$\omega^*(j)' = \left(\omega^*(j, \frac{1}{m}) \quad \omega^*(j, \frac{2}{m}) \quad \dots \quad \omega^*(j, 1) \right)$$

and normalize the vector by dividing each element by $\omega^*(j)'\omega^*(j)$ (i.e., the column sum). This yields the normalized weights vectors $\omega(0)$ and $\omega(1)$.

(iii) Stack the weights across states of E and G :

- By treatment state j :

$$\omega' = \begin{pmatrix} \omega(0)' & \omega(1)' \end{pmatrix} \in \mathbb{R}^{1 \times 2m}$$

- By genetic stratum g :

$$\omega^{G'} = \begin{pmatrix} -\omega' & \mathbf{0}_m & \mathbf{0}_m & \mathbf{0}_m & \omega' \end{pmatrix} \in \mathbb{R}^{1 \times 10m}$$

The weights take the difference in the effects of E averaged over $u \in \{0.55, 0.85\}$ between the last and the first column, that is, $g = 5$ and $g = 1$.

Remark: The shape restriction C. “Additive separability” ensures that effects are also identified for $G \in \{2, 3, 4\}$ —we document that these quintiles are reliably estimated in our robustness checks.

4. Set target function (a scalar):

$$\left(\omega^{G'}(\theta' b^{E \times G})\right)' \in \mathbb{R}^{1 \times 1}$$

Remark: The vector $\theta \in \mathbb{R}^{10n \times 1}$ contains all $10n$ parameters that are computed through linear optimization, and therefore unspecified in this stage.

The vector $\theta' b^{E \times G} \in \mathbb{R}^{10m \times 1}$ gives the 10 stacked MTRs (evaluated at the elements of u) for (for combinations of G and E), which $\omega^{G'}$ aggregates to the difference in average effects of E between the last and first quintile of G —reducing the dimensionality of this problem to a scalar value.

5. Set constraints:

- (i) Support constraints:

$$\theta' b^{E \times G} \leq 20 \quad \text{and} \quad \theta' b^{E \times G} \geq 0$$

(ii) Moment constraints (MTRs must reproduce moments):

$$\mathbf{w}^{IV'}(\boldsymbol{\theta}'\mathbf{b}^{E \times G}) = \boldsymbol{\delta},$$

where the weighting matrix $\mathbf{w}^{IV} \in \mathbb{R}^{10m \times 20}$ for the 20 moments (all combinations between genetic stratum g , treatment state j and instrument value z , i.e., $5 \times 2 \times 2 = 20$) is constructed as follows:

- First, we construct dummy variables that indicate whether a group with the observed values g , z , and j may stem from a specific value u :

$$w(u, z, 1, g)^* = \mathbb{1}\{u \leq \pi_{g,z}\}$$

$$w(u, z, 0, g)^* = \mathbb{1}\{u \geq (1 - \pi_{g,z})\}$$

- Then, we assemble all elements in column vectors within treatment state j , genetic stratum g and the instrument value z , yielding $\mathbf{w}(z, j, g)^*$, and normalize by dividing by their treatment-state specific sum $(\mathbf{w}(z, j, g)^*{}'\mathbf{w}(z, j, g)^*)$. This yields $\mathbf{w}(z, j, g)$.
- For each j and g , we construct the matrix

$$\mathbf{w}(j, g) = \begin{pmatrix} \mathbf{w}(0, j, g) & \mathbf{w}(1, j, g) \end{pmatrix} \in \mathbb{R}^{m \times 2}.$$

- We then assemble these 10 different matrices in a block-diagonal matrix as follows:

$$\mathbf{w}^{IV} = \text{bdiag} \left(\mathbf{w}(0, 1), \mathbf{w}(1, 1), \mathbf{w}(0, 2), \mathbf{w}(1, 2), \dots, \mathbf{w}(0, 5), \mathbf{w}(1, 5) \right) \in \mathbb{R}^{10m \times 20}$$

(iii) Shape constraints:

For all constraints, we construct a slope matrix $\mathbf{S} \in \mathbb{R}^{m \times m-1}$, with the following elements in row k and column l :

$$S_{kl} = \begin{cases} -1 & \text{if } k = l \\ 1, & \text{if } k + 1 = l \\ 0, & \text{otherwise.} \end{cases}$$

Remark: This matrix is valuable for a single MTR. For each MTR, it adds $m - 1$ restriction, comparing the level of an MTR at the point u_i with the one at $u_{i-1} \quad \forall 2 \leq i \leq m$.

(iv) We then use this matrix and specify the following shape constraints:

A. Monotone treatment selection for $E = 1$:

$$\mathbf{A}_1 = \mathbf{I}(5) \otimes \begin{pmatrix} 0 \\ 1 \end{pmatrix} \otimes \mathbf{S}$$

Remark: The first term $\mathbf{I}(5)$ applies the restriction to each genetic stratum, the second applies it to treatment state $E = 1$, but switches

it off for treatment state $E = 0$, the third is the actually the slope restriction.

B. No selection into losses:

$$\mathbf{A}_2 = \mathbf{I}(5) \otimes \begin{pmatrix} -1 \\ 1 \end{pmatrix} \otimes \mathbf{S}$$

Remark: The difference to A. is the second term, ensuring that, within G , the slope restrictions are applied to the difference between the treatment states.

C. Additive separability:

- Generate a matrix that restricts the slope between cells of g : $\mathcal{G} \in \mathbb{R}^{5 \times 5-1}$, with its elements in row k and column l defined as follows:

$$\mathcal{G}_{kl} = \begin{cases} 1 & \text{if } k = 1 \\ -1, & \text{if } k = l + 1 \\ 0, & \text{otherwise.} \end{cases}$$

- With this matrix, the restriction reads:

$$\mathbf{A}_3 = \mathcal{G} \otimes \mathbf{I}(2) \otimes \mathbf{S}$$

Remark: The difference to A. is the first term, ensuring that the slope restrictions are also applied between G .

D. Overall conditions that we can use as equality restrictions read:

$$wA_R = \begin{pmatrix} (\theta'b)' A_1 \\ (\theta'b)' A_2 \\ (\theta'b)' A_3 \end{pmatrix}$$

However, to be more flexible in the optimization process (see step B-7.), we use the following conditions as inequality restrictions

$$wA_R = \begin{pmatrix} (\theta'b)' A_1 \\ -(\theta'b)' A_1 \\ (\theta'b)' A_2 \\ -(\theta'b)' A_2 \\ (\theta'b)' A_3 \\ -(\theta'b)' A_3 \end{pmatrix}$$

6. We optimize the target function (i.e., minimize or maximize) subject to the constraints by choosing θ :

$$\begin{aligned}
 & \max_{\theta} / \min_{\theta} \omega^{G' \theta' b} \\
 & \text{s.t.} \\
 & \quad w' \theta' b = \delta \quad (\text{equality constraints}) \\
 & \quad \begin{pmatrix} \theta' b \\ -\theta' b \end{pmatrix} \leq \begin{pmatrix} 20 \\ \vdots \\ 20 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \\
 & \quad w A_R \leq f \quad (\text{inequality constraints})
 \end{aligned}$$

To ensure that the last set of inequality restrictions hold with equality, we can set $f = \mathbf{0}$ —the preferred value.

Finally, we use `mata's LinearProgram()` command in Stata to solve this problem.

7. We start solving the problem with $f = \mathbf{0}$. If the problem is infeasible, we start allowing small deviations from 0 for the restriction C and increase the elements in f for the respective restrictions until the problem can be solved.

Remark: This is similar to [Rose and Shem-Tov \(2021\)](#). In contrast to their approach, we do not first set up an auxiliary problem to determine the optimal values of \mathcal{f} .

Part C: Bootstrap

1. For inference, execute Part A B times with resampled data. This yields B different π^b and δ^b vectors, indexed by the bootstrap iteration $b \in \{1, \dots, B\}$.
2. With this input, we execute Part B using the resampled moments (including the shares), get θ_{\max}^b , θ_{\min}^b , compute the corresponding MTRs and g -specific treatment effects. Finally, we use the standard deviation of the corresponding estimate as the standard error for statistical inference.

Figure [F.3](#) plots the results of this approach, where the slightly transparent, horizontal lines are the "moments" (G_i -specific outcome means and their placement on the unit-interval, which we derive from the complier shares). The blue (for the treated outcome) and red (for the untreated) lines are the output of this linear programming approach. They reflect the minimal (the dashed lines) and maximal (the solid lines) possible interaction effect (defined in the main text) that the MTR lines produce while being consistent with the shape restrictions and matching the moments.

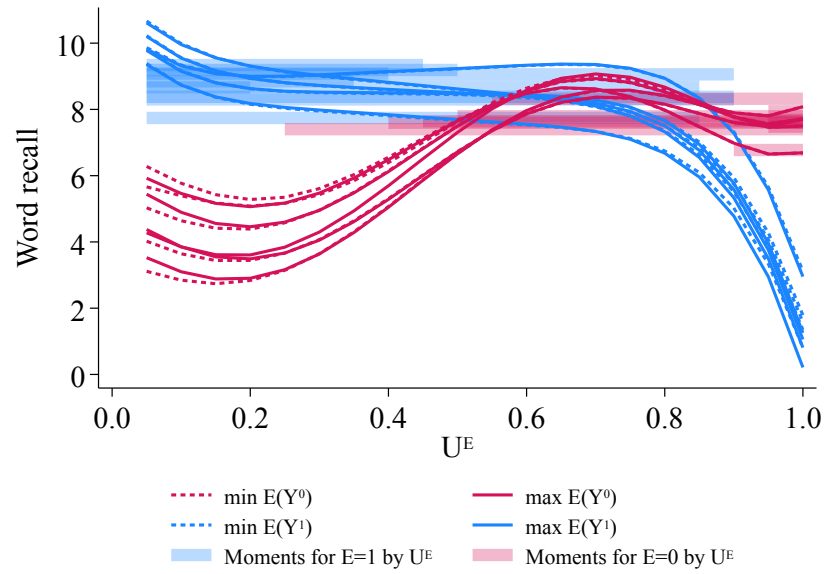


Figure F.3: Potential outcome curves estimated with Bernstein polynomials

Notes: This figure shows the minima and maxima of the ten potential outcome curves estimated via linear program with Bernstein polynomials using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. Blue indicates curves and moments for $E_i = 1$, and red indicates $E_i = 0$. Solid lines are maxima; dashed lines are minima of the potential outcome curves. There are five pairs of curves for $E_i = 1$ and five for $E_i = 0$, one pair for every PGI quintile. Every pair consists of a minimum and a maximum that bound the potential outcome curve for its respective quintile. The vertical bars indicate the moments the curves must pass and the U^E ranges of individuals contributing to these means.

G Testing “no selection into losses” (non-positive MTE slopes)

A critical constraint we apply in our linear programming approach is “no selection-into-losses”, i.e., no MTEs that increase in U^E . To test this in our setting, we follow [Imbens and Rubin \(1997\)](#) and use the instrument to compute mean outcomes for always-takers, treated and untreated compliers, and never-takers. For simplicity, we test this condition globally and do not distinguish between cells of G_i (we show the complete G_i -specific means in [Figure F.1](#)). We present the results in [Table G.1](#). In Panel A, we focus on differences between always-takers and treated compliers (Column 3) and untreated compliers and never-takers (Column 6). The differences are informative about whether the treated outcome $\mathbb{E}[Y_i^1|U^E = u]$ and the untreated outcome $\mathbb{E}[Y_i^0|U^E = u]$ —the difference of which is the MTE—are heterogeneous in U^E .

Column (3) presents the mean word recall difference between always-takers and treated compliers. It shows a substantial and statistically significant heterogeneity: Always-takers recall about 1.25 words more. Intuitively, this is unsurprising, as always-takers to a compulsory schooling reform will, on average, have more years of education, will be more likely to hold advanced degrees, or may be positively selected in terms of unobserved characteristics (if we have selection into gains, which we want to argue). Furthermore, this result shows that $\mathbb{E}[Y_i^1|U^E = u]$ has a negative slope. Likewise, we do the same with untreated compliers and never-takers. Here, the heterogeneity is less pronounced and not statistically significant. If we conclude that both groups do not have different outcomes, we can stop as in this case, the

difference in the first two groups proves that we have essential heterogeneity. If the insignificant difference is meaningful, things may change. The difference is also negative, contrasting the existing empirical evidence for the slope of the untreated outcome (see, e.g., [Carneiro and Lee, 2009](#); [Westphal et al., 2022](#)). However, it is essential to mention that never-takers should not exist with a compulsory schooling reform, where everyone should be forced to stay in school until age 15. If this group has never existed, this might be a measurement error. If these individuals had special exemptions from the rule change (and therefore existed), the difference between never-takers and untreated compliers may not inform about the global course of the curve. Assessing the multiple complier groups that we gain by stratifying by G_i (see Figure F.1) indeed suggests that never-takers are different and $\mathbb{E}[Y_i^0|U^E = u]$ indeed increases when $U^E < 0.95$.

Nonetheless, with only a binary instrument and without exploiting covariate heterogeneity together with the additive separability assumption (which we will do below), an additional linearity assumption is necessary (due to the never-takers) to point-identify a marginal treatment effect via the method introduced by [Brinch et al. \(2017\)](#). We document a formal test of the slope of $\mathbb{E}[Y_i^1|U^E = u]$ and $\mathbb{E}[Y_i^1 - Y_i^0|U^E = u]$ in Panel B.²³ It shows that the slope of the treated outcome is negative and statistically significant (as shown in Panel A). The slope of the linear MTE is also negative and still large in magnitude. However, likely due to the concerns about never-takers outlined above, it is not statistically significant, albeit with a negative sign. Again, evidence from the G_i -specific complier groups strongly suggests that

²³The exact formula reads

$$\frac{\partial \mathbb{E}[Y_i^1|U^E = u]}{\partial U^E} = \frac{Y_i^{CT} - Y_i^{AT}}{\frac{\pi^C + \pi^{AT}}{2}}, \quad \frac{\partial \mathbb{E}[Y_i^1 - Y_i^0|U^E = u]}{\partial U^E} = \frac{Y_i^{CT} - Y_i^{AT}}{\frac{\pi^C + \pi^{AT}}{2}} - \frac{Y_i^{NT} - Y_i^{CU}}{\frac{\pi^C + \pi^{NT}}{2}},$$

where Y_i^{AT} , Y_i^{CT} , Y_i^{CU} , and Y_i^{NT} are means from Columns (1), (2), (4), and (5), respectively and π^{AT} , π^C , π^{NT} are the corresponding shares (compliers do not need to be differentiated).

Table G.1: Mean outcomes by instrument response types and test for essential heterogeneity

	Unobserved heterogeneity					
	in the treated outcome			in the untreated outcome		
	(1) Always-takers	(2) Treated compliers	(3) Difference (2) – (1)	(4) Untreated compliers	(5) Never-takers	(6) Difference (5) – (4)
<i>Panel A:</i>						
Mean word recall:	9.500 (0.215)	8.306 (0.332)	–1.245*** (0.454)	8.109 (0.215)	7.679 (0.340)	–0.353 (0.396)
Share:	0.456 (0.035)	0.489 (0.036)		0.489 (0.036)	0.055 (0.011)	
<i>Panel B:</i>						
Test for essential heterogeneity: (sufficient condition, may be uninformative if heterogeneity is nonlinear)						
Slope of $\mathbb{E}[Y_i^1 U^E = u]$			–2.631*** (0.961)			
Slope of MTE $\mathbb{E}[Y_i^1 - Y_i^0 U^E = u]$			–1.326 (1.423)			

Notes: This table presents estimates of mean outcomes for always-takers, treated and untreated compliers, and never-takers (panel A) as well as results of a test for essential heterogeneity (panel B) using data from ELSA waves 1–6 and our main sample selection, as outlined in Section 2.2. We compute the type-specific shares using the specification of Eq. (2) without G_i . The complier share is the coefficient on Z_i , the always-taker share is the constant (as all variables are demeaned), and the never-taker share is the remainder. For the type-specific outcome means, we compute means by E_i and Z_i (and their interaction) using a reduced-form specification to regress word recall on the same controls and full interactions of E_i and Z_i . As compliers never appear alone in these means, we use the formula provided in Imbens and Rubin (1997) and the type-specific shares. Standard errors are computed using 200 bootstrap replications and are shown in parentheses. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$ indicate significance levels for the differences.

the $\mathbb{E}[Y_i^0|U^E = u]$ increases at least for a relevant range when $U^E < 0.95$. We conclude that we likely face essential heterogeneity in our setting. Combined differences in the first stage induced by G_i , the result may be that 2SLS cannot recover the true interaction parameter we would need to make accurate statements about the interaction effect.