






# Sibling comparisons elucidate the associations between educational attainment polygenic scores and alcohol, nicotine and cannabis

Jessica E. Salvatore<sup>1,2</sup> , Peter B. Barr<sup>1</sup> , Mallory Stephenson<sup>1</sup> , Fazil Aliev<sup>1,3</sup>, Sally I-Chun Kuo<sup>1</sup>, Jinni Su<sup>4</sup>, Arpana Agrawal<sup>5</sup>, Laura Almasy<sup>6,7</sup>, Laura Bierut<sup>5</sup>, Kathleen Bucholz<sup>5</sup>, Grace Chan<sup>8</sup>, Howard J. Edenberg<sup>9</sup>, Emma C. Johnson<sup>5</sup> , Vivia V. McCutcheon<sup>5</sup> , Jacquelyn L. Meyers<sup>10</sup>, Marc Schuckit<sup>11</sup>, Jay Tischfield<sup>12</sup>, Leah Wetherill<sup>13</sup> & Danielle M. Dick<sup>1,14,15</sup>

## ABSTRACT

**Background and Aims** The associations between low educational attainment and substance use disorders (SUDs) may be related to a common genetic vulnerability. We aimed to elucidate the associations between polygenic scores for educational attainment and clinical criterion counts for three SUDs (alcohol, nicotine and cannabis). **Design** Polygenic association and sibling comparison methods. The latter strengthens inferences in observational research by controlling for confounding factors that differ between families. **Setting** Six sites in the United States. **Participants** European ancestry participants aged 25 years and older from the Collaborative Study on the Genetics of Alcoholism (COGA). Polygenic association analyses included 5582 (54% female) participants. Sibling comparisons included 3098 (52% female) participants from 1226 sibling groups nested within the overall sample. **Measurements** Outcomes included criterion counts for DSM-5 alcohol use disorder (AUDSX), Fagerström nicotine dependence (NDSX) and DSM-5 cannabis use disorder (CUDSX). We derived polygenic scores for educational attainment (*EduYears-GPS*) using summary statistics from a large (> 1 million) genome-wide association study of educational attainment. **Findings** In polygenic association analyses, higher *EduYears-GPS* predicted lower AUDSX, NDSX and CUDSX [ $P < 0.01$ , effect sizes ( $R^2$ ) ranging from 0.30 to 1.84%]. These effects were robust in sibling comparisons, where sibling differences in *EduYears-GPS* predicted all three SUDs ( $P < 0.05$ ,  $R^2$  0.13–0.20%). **Conclusions** Individuals who carry more alleles associated with educational attainment tend to meet fewer clinical criteria for alcohol, nicotine and cannabis use disorders, and these effects are robust to rigorous controls for potentially confounding factors that differ between families (e.g. socio-economic status, urban–rural residency and parental education).

**Keywords** Alcohol, cannabis, Collaborative Study on the Genetics of Alcoholism, nicotine, polygenic risk score, sibling comparisons.

Correspondence to: Jessica E. Salvatore, Department of Psychology and the Virginia Institute for Psychiatric and Behavioral Genetics, VCU PO Box 842018, 806 West Franklin Street, Richmond, VA 23284-2018, USA. E-mail: jesalvatore@vcu.edu

Submitted 27 February 2019; initial review completed 30 May 2019; final version accepted 2 September 2019

## INTRODUCTION

Researchers have studied the associations between educational attainment and substance use disorders (SUDs) for more than a century [1,2]. Cross-sectional studies consistently link use of tobacco, alcohol and cannabis with high school dropout [2] and greater educational attainment with lower rates of SUD diagnoses [3–5]. There is a substantial body of work exploring the hypotheses that SUDs influence early termination of education and that early termination of education influences SUDs [6], with evidence supporting both temporal orderings [5,7–10]. A third hypothesis is also plausible: that the

associations between low educational attainment and SUDs are attributable, at least in part, to a common general vulnerability. Genetic factors represent one type of general vulnerability. Consistent with this possibility, genetic epidemiological data indicate that there is a set of genetic factors that influence both low educational attainment and a higher likelihood of developing SUDs [11–14]. There is also evidence that familial factors confound the associations between educational attainment and multiple forms of substance use and dependence [6,15], although the specific source of this familial confounding (i.e. genes or the rearing environment) was not specified in those studies.

Recent advances in characterizing the molecular genetic basis of complex traits and behaviors have stimulated interest in translating findings from genetic epidemiological studies, which use patterns of resemblance among individuals of known genetic relatedness to make inferences about latent genetic influences on traits and behaviors, into a molecular genetic framework [16,17]. This is typically accomplished using a polygenic scoring approach, where researchers leverage genome-wide association results from large, well-powered discovery samples to calculate personalized indices of the weighted number of trait-associated alleles carried by each participant in an independent sample [18,19]. In polygenic analyses, the associations between these polygenic scores and other traits and behaviors are examined to determine their shared genetic etiology.

In this study, we combined polygenic association and sibling comparison methods to elucidate the associations between polygenic scores for educational attainment [20] and clinical criterion counts for three common SUDs (alcohol, nicotine and cannabis) in a sample of adults of European ancestry. Sibling comparisons [21–24] provide a complementary tool to clarify the nature of associations observed in polygenic analyses. Biological full siblings reared together share the same home environment and a substantial portion of their genetic variation (50% on average), allowing for control of measured and unmeasured familial factors such as socio-economic status, religious upbringing, urban–rural residency, parental education and familial polygenic load, that are also known to influence SUD outcomes. Controlling for these potential confounders shared by siblings is important because too often polygenic associations are over-interpreted as evidence that a particular set of alleles has pleiotropic effects across traits or disorders. For this reason, testing the alternative explanation that polygenic associations are attributable to familial confounding is important for understanding the molecular genetic basis underlying the links between low educational attainment and SUDs. This is particularly critical in view of the enthusiasm to incorporate polygenic scores as part of precision medicine efforts to identify and intervene with individuals deemed genetically at risk.

Significant associations between an educational attainment polygenic score and SUD criterion counts within a sibling comparison design would be consistent with the interpretation that carrying more alleles associated with educational attainment is associated with a lower likelihood of developing SUD problems. In contrast, if sibling differences in educational attainment polygenic scores do not predict SUD criterion counts it suggests that polygenic associations are confounded by other shared familial factors. This difference is important, considering that social advantage is related to both

educational attainment polygenic scores [25–27] and rates of SUDs [28].

## MATERIALS AND METHODS

### Participants

Participants came from the Collaborative Study on the Genetics of Alcoholism (COGA) [29–31], whose objective is to identify genes involved in alcohol dependence and related disorders. Proband (i.e. index individuals) were identified through alcohol treatment programs at six US sites. Probands and their families were invited to participate if the family was sufficiently large (usually sibships > 3 with parents available), with two or more members in the COGA catchment area. Comparison families were recruited from the same communities. The Institutional Review Boards at all data-gathering sites approved this study and written consent was obtained from all participants. COGA data are available via dbGaP (phs000763.v1.p1, phs000125.v1.p1) or through the National Institute on Alcohol Abuse and Alcoholism.

We defined two study samples within COGA. The first sample included all participants of European ancestry aged 25 years or older with both genome-wide association data and relevant SUD phenotypic information [ $n = 5582$  individuals from 1093 extended families; 3009 (54%) female;  $\text{mean}_{\text{age}} = 42.29$  years, age range = 25–91 years]. We limited the sample to those of European ancestry to avoid population stratification [32] because the educational attainment genome-wide association study (GWAS) weights come from a European ancestry discovery sample. SNPrelate [33] was implemented to estimate principal components from GWAS data and subsequently used to determine European ancestry. We implemented the age minimum to balance the needs to ensure that the majority of participants had passed through the period of highest risk for onset of the SUDs without unduly limiting sample size. Epidemiological data regarding age of onset for SUDs [34–36] guided our decision to select age 25 as the cut-off, which also mirrors the cut-off used in analyses of educational attainment in US Census data [37].

The second sample was a subset of the first sample, limited to groups of European-ancestry biological full siblings (confirmed by genotyping) nested within the larger COGA sample. This process identified 4733 individuals nested within 1655 sibling groups (two to 12 siblings per group). As detailed below, the 4733 sibling GWAS samples were used to calculate the educational attainment polygenic scores used for the sibling comparison analyses. The sample was subsequently filtered by age at phenotypic assessment for the linear mixed-model analyses. In total, the sibling comparison analyses included 3098 individuals [1616 (52%) female] who were 25 years of age or older

( $\text{mean}_{\text{age}} = 37.89$  years) from 1226 sibling groups nested within 773 extended families.

## Measures

### Genotyping

Genotyping for the COGA European ancestry participants was performed using the Illumina 1 M, Illumina OmniExpress (Illumina, San Diego, CA, USA) and Smokescreen (BioRealm, Walnut, CA, USA) arrays. Quality control and imputation procedures are described in Lai *et al.* [31] and in the Supporting information, section S1.

### Substance use disorder clinical criterion counts

Clinical criterion counts for alcohol (AUDSX), nicotine (NDSX) and cannabis (CUDSX) were obtained from the reliable and valid Semi-Structured Assessment for the Genetics of Alcoholism (SSAGA) [38,39]. Criterion counts for alcohol and cannabis use disorder were made according to DSM-5 [40], and thus each had a possible range of 0–11. The criterion count for NDSX came from the Fagerström Test for Nicotine Dependence [41] and had a possible range of 0–10. The criterion count distributions showed right skewness; to address this in inferential analyses, we applied a logarithmic transformation (left anchored at 1).

### Covariates and measures for robustness and sensitivity analyses

We included sex, age at last interview, cohort [indexed using three dummy-coded variables derived from participant year of birth: (1896–1930) set as reference; (1930–50); (1950–70); (1970–2010)] and the first two principal components for genetic ancestry in all analyses.

We conducted a series of robustness and sensitivity analyses to probe and interpret the effects from our primary analyses. For robustness analyses, we used participants' educational attainment, assessed as highest level of education completed. Potential responses ranged from 0 to 17 years (primary or secondary school = actual year; technical school/1 year college = 13 years; 2 years college = 14 years; 3 years college = 15 years; 4 years college = 16 years; any graduate degree = 17 years). In sensitivity analyses of the sibling data, we used participants' reports of their living arrangements while growing up from a set of 13 options (see Supporting information, section S2) to evaluate whether the pattern of effects changed when the sample was limited to siblings who reported the same living arrangements (and thus probably shared the same rearing environment). An early version of the SSAGA did not query living arrangements; accordingly, we were only able to confirm that siblings grew up

together for a subset of the sample (see Supporting information, section S2).

## STATISTICAL METHODS

### Educational attainment genome-wide polygenic scores (*EduYears-GPS*)

We used results from the Social Science Genetic Association Consortium (SSGAC) GWAS of educational attainment [20] to construct educational attainment genome-wide polygenic scores (*EduYears-GPS*) in the COGA sample. Although polygenic scores are often described as polygenic risk scores, we prefer the term 'genome-wide polygenic score' for this study. This is because 'risk' connotes a negative outcome, whereas educational attainment is typically valued. After removing palindromic single nucleotide polymorphisms (SNPs) (which can be ambiguous with respect to the reference allele in different samples), we used the *clump* and *score* procedures in PLINK [42] to sum each individual's total number of minor alleles from the score SNPs, with each SNP weighted by the negative log of the GWAS association *P*-value and sign of the association (beta) statistic. Clumping was performed with respect to the linkage disequilibrium (LD) pattern in the COGA EA sample (founders only) using a 500-kb physical distance and an LD threshold of  $r^2 \geq 0.25$ . Following conventions for polygenic scoring using the pruning-and-thresholding approach [18], we calculated a series of GPS in COGA that included SNPs meeting increasingly stringent *P*-value thresholds in the discovery GWAS ( $P < 0.50$ ,  $P < 0.40$ ,  $P < 0.30$ ,  $P < 0.20$ ,  $P < 0.10$ ,  $P < 0.01$ ,  $P < 0.001$ ,  $P < 0.0001$ ).

### Association of *EduYears-GPS* and SUDS

We examined associations between *EduYears-GPS* and the SUD criterion counts in separate linear mixed models using the *nlme* package version 3.1–128 [43] for R version 3.2.3 [44]. We conducted preliminary analyses to identify the *EduYears-GPS* most strongly associated with criterion counts for each SUD (see Supporting information, Table S1), and present results using the threshold with the strongest association. We conducted these preliminary analyses separately for polygenic scores meeting increasingly stringent *P*-value thresholds using linear mixed models, which allowed us to account for the nested structure of the COGA family-based data; other methods for optimizing the *P*-value threshold, e.g. PRSice [45], do not allow for nested data. In addition to the covariates described above, we also included a count measure of the number of SNPs available for scoring for each participant. Marginal effect sizes for fixed effects were calculated using the MuMIn package version 1.15.6 [46].

### Sibling comparisons of *EduYears-GPS* and SUDs

We used the 4733 sibling GWAS sample to calculate the *EduYears-GPS-mean* (for each sibling group) and *EduYears-GPS-deviation* scores (for each individual within that sibling group). We then filtered the sample based on participants' age at last interview to retain those who were aged 25 years or older (age cut-off selected to ensure that participants had passed through the period of highest risk for onset of the SUDs) for our primary sibling comparisons sample; additional information regarding this process can be found in Supporting information, section S3. Using all available GWAS data from a sibling group to calculate the *EduYears-GPS-mean* and *EduYears-GPS-deviation* scores has the advantage of providing a more precise estimate for these variables (as genotype does not change with age) versus limiting calculation of *EduYears-GPS-mean* to those siblings who also met the phenotypic age threshold. In separate linear mixed models, we then examined whether *EduYears-GPS-deviation* predicted SUDs after controlling for *EduYears-GPS-mean*. The sibling comparison is captured by the *EduYears-GPS-deviation* parameter, and indicates whether sibling differences in *EduYears-GPS* predict SUDs; this parameter captures the within-family effect. The *EduYears-GPS-mean* parameter captures whether family-level differences in *EduYears-GPS* predict SUDs, reflecting the between-family effect.

### Robustness and sensitivity analyses

We conducted robustness analyses to examine whether findings changed when statistically controlling for educational attainment in both the association and sibling comparison analyses. Sibling differences and family means for phenotypic educational attainment (i.e. *EduYears-deviation* and *EduYears-mean*) were calculated using the same procedure described above for *EduYears-GPS-deviation* and *EduYears-GPS-mean*.

We conducted sensitivity analyses to see whether effects changed when using a more conservatively defined subsample of siblings who were known to have the same living arrangements while growing up or who were born within 3 years of the eldest. These more conservative definitions assume that siblings who report the same living arrangements growing up and who are born in closer proximity to one another are likely to share more features of their home environment than siblings who report different living arrangements or who are born further apart. In total, 1702 individuals (54% female) from 739 sibling groups were available for this analysis. We also examined whether the effects were robust when sibships that included monozygotic twins (eight sibling groups) were removed from the analysis. Monozygotic twins share 100% of their genetic variation, and we wanted to ensure that our results were not driven by genotyping errors or

PLINK's handling of SNPs set to missing (as part of cleaning for Mendelian errors) during polygenic score calculation. Sample size as a function of the filters employed for these sensitivity analyses are shown in Supporting information, Fig. S1.

## RESULTS

### Descriptive statistics

Descriptive statistics for the SUD criterion counts and educational attainment for the full sample ( $n = 5582$ ) and the sibling subsample ( $n = 3098$ ) are summarized in Table 1. Representativeness analyses of the sibling subsample are summarized in Supporting information, section S4.

### Polygenic association for *EduYears-GPS* and SUDs

We identified the  $P < 0.30$  threshold for AUDSX,  $P < 0.20$  for NDSX and  $P < 0.01$  for CUDSX as the *EduYears-GPS* thresholds most strongly associated with each SUD criterion count. As shown in Table 2, higher *EduYears-GPS* was associated with lower SUD criteria. The *EduYears-GPS* accounted for 0.79, 1.84 and 0.30% of the variance in AUDSX, NDSX and CUDSX, respectively.

### Sibling comparisons of *EduYears-GPS* and SUDs

We carried forward the substance-specific thresholds that were most strongly associated with each criterion count from above into the sibling comparisons to examine

**Table 1** Descriptive statistics.

<i>Full sample (n = 5582; 54% female)</i>					
<i>Measure</i>	<i>n</i>	<i>Mean</i>	<i>SD</i>	<i>Min</i>	<i>Max</i>
Age (years)	5582	42.29	13.26	25	91
AUDSX	5582	3.77	3.83	0	11
NDSX	4754	2.61	3.00	0	10
CUDSX	5578	1.56	2.81	0	11
Educational Attainment (years)	5578	13.43	2.33	2	17
<i>Sibling subsample (n = 3098; 52% female)</i>					
<i>Measure</i>	<i>n</i>	<i>Mean</i>	<i>SD</i>	<i>Min</i>	<i>Max</i>
Age (years)	3098	37.89	10.85	25	81
AUDSX	3098	4.39	3.90	0	11
NDSX	2752	2.58	3.00	0	10
CUDSX	3097	1.94	3.04	0	11
Educational Attainment (years)	3095	13.55	2.29	5	17

SD = standard deviation; AUDSX = DSM-5 alcohol use disorder; NDSX = Fagerström nicotine dependence; CUDSX = DSM-5 cannabis use disorder.

**Table 2** Associations between *EduYears-GPS* genome-wide polygenic scores and substance use disorder criterion counts in the full sample.

Parameter	Alcohol use disorder criterion count ( <i>n</i> = 5582)		Nicotine dependence criterion count ( <i>n</i> = 4754)		Cannabis use disorder criterion count ( <i>n</i> = 5578)	
	<i>b</i>	95% CI	<i>b</i>	95% CI	<i>b</i>	95% CI
Intercept	-2.79	(-3.78, -1.79)	-1.47	(-2.59, -0.34)	-1.01	(-1.74, -0.28)
Sex (female)	-0.58	(-0.62, -0.54)	-0.22	(-0.27, -0.18)	-0.34	(-0.37, -0.30)
Age (years)	-0.01	(-0.01, -5.12E-03)	<b>5.00E-03</b>	(-3.84E-05, 9.76E-03)	<b>-8.60E-03</b>	(-0.01, -5.01E-03)
PC1	<b>83.39</b>	(9.18, 157.60)	-23.63	(-107.81, 60.54)	-42.38	(-108.61, 23.86)
PC2	26.39	(-11.30, 64.07)	29.71	(-12.40, 71.83)	15.99	(-17.76, 49.73)
<i>EduYears-GPS</i> count	<b>1.01E-05</b>	(8.35E-06, 1.19E-05)	<b>9.29E-06</b>	(6.78E-06, 1.18E-05)	<b>3.70E-05</b>	(2.77E-05, 4.62E-05)
Cohort 2	0.22	(0.10, 0.33)	0.15	(-0.01, 0.31)	-0.03	(-0.14, 0.07)
Cohort 3	<b>0.41</b>	(0.25, 0.58)	0.14	(-0.08, 0.35)	<b>0.44</b>	(0.29, 0.58)
Cohort 4	0.16	(-0.05, 0.36)	-0.02	(-0.28, 0.24)	<b>0.31</b>	(0.12, 0.49)
<i>EduYears-GPS</i>	<b>-19360.64<sup>a</sup></b>	(-25072.31, -13648.97)	<b>-24663.58<sup>a</sup></b>	(-29961.23, -19365.93)	<b>-2551.10</b>	(-3781.71, -1320.49)
<i>EduYears-GPS</i> $\Delta R^2$	<b>0.79%</b>		<b>1.84%</b>		<b>0.30%</b>	

Bold type indicates estimate  $P \leq 0.05$ . <sup>a</sup>The *EduYears-GPS* effect was robust after controlling for phenotypic educational attainment (see Supporting information, Table S2). The *EduYears-GPS* thresholds for each substance were: alcohol ( $P < 0.30$ ); nicotine ( $P < 0.20$ ); cannabis ( $P < 0.01$ ). PC = principal component for genetic ancestry; cohort = dummy-coded variables indexing year of birth, defined as (1930–50); (1950–70); and (1970–2010); *EduYears-GPS* count = number of single nucleotide polymorphisms available for polygenic scoring; *EduYears-GPS* = educational attainment genome-wide polygenic score;  $\Delta R^2$  = change in  $r$ -squared; CI = confidence interval.

whether *EduYears-GPS*-deviation predicted each SUD criterion count after controlling for *EduYears-GPS*-mean.

The results of the sibling comparisons are shown in Table 3. Individuals with higher *EduYears-GPS* compared to their siblings had lower alcohol, nicotine and cannabis criterion counts. Sibling differences in *EduYears-GPS* accounted for 0.17, 0.20 and 0.13% of the variance in AUDSX, NDSX and CUDSX, respectively. There were also family-level effects, whereby those in sibling groups with higher *EduYears-GPS*-mean had lower alcohol, nicotine and cannabis criterion counts. These family-level effects accounted for 0.29, 1.89 and 0.22% of the variance in AUDSX, NDSX and CUDSX, respectively.

### Robustness analyses

After controlling for participants' measured (phenotypic) educational attainment in the polygenic analyses, *EduYears-GPS* continued to be associated with AUDSX and NDSX (but not CUDSX) (Supporting information, Table S2). After controlling for sibling and family differences in educational attainment in the sibling comparison analyses, the effects of sibling differences in *EduYears-GPS* on SUD criterion counts were attenuated for NDSX and CUDSX ( $P = 0.09$ – $0.13$ ), but remained significant for AUDSX ( $P = 0.01$ ) (Supporting information, Table S3). Sibling and family differences in educational attainment were also significantly associated with SUD criterion counts. Individuals with higher educational attainment compared to their siblings and individuals from sibling groups with higher educational attainment had lower AUDSX, NDSX and CUDSX.

### Sensitivity analyses

In the first set of sensitivity analyses, we examined whether effects held when the sample was limited to the groups of siblings who were known to have grown up together ( $n = 739$  sibling groups). In the second set of sensitivity analyses, we examined whether the effects held when the sample was limited to those who were born within 3 years of the first-born in a sibling group. In the third set of sensitivity analyses, we examined whether the effects were also robust when sibships that included monozygotic twins (eight sibling groups) were removed from the analysis. Across all three sets of sensitivity analyses in smaller, more conservative test samples, we continued to find that individuals with higher *EduYears-GPS* than their siblings had lower SUD criterion counts (Supporting information, Tables S4–S6). The only exception to this was that the effect of sibling differences in *EduYears-GPS* on CUDSX was attenuated ( $P = 0.08$ ) in the sensitivity analyses limited to those born within 3 years of the first born in a sibling group.

## DISCUSSION

The present study illustrates how sibling comparisons can improve our understanding of the shared genetic etiology underlying educational attainment and substance use problems. Consistent with previous findings that educational attainment has a negative genetic correlation with alcohol problems [11,13], cannabis use disorder [14] and smoking [12], we found that individuals met fewer SUD criteria when they carried more alleles associated with educational attainment. We replicated these effects within a sibling comparisons design, where we found that individuals met fewer clinically significant substance use criteria when they carried more alleles associated with higher educational attainment than their siblings. Sibling comparisons are uniquely powerful because they control for unmeasured confounding factors shared by siblings that could otherwise explain the association between educational attainment polygenic scores and substance use disorder criteria: factors such as socio-economic status, urban–rural residency and parental education. Thus, our findings suggest that the association between educational attainment polygenic scores and SUDs is not completely explained by confounders that differ between families.

These findings add important nuance to discussions regarding the nature of associations between educational attainment and problematic substance use. First, our findings are consistent with previous findings that educational outcomes reflect many genetically influenced traits and behaviors, including SUD-associated factors such as behavior problems, attention-deficit hyperactivity disorder, and personality [25,26,47–50], not simply intelligence or cognitive ability. Interestingly, in our robustness analyses, the educational attainment polygenic scores predicted alcohol use disorder and nicotine dependence criterion counts above and beyond participants' observed (phenotypic) educational attainment. This highlights that these polygenic scores index factors linked to educational persistence and SUDs that are not fully captured by educational attainment itself. In contrast, for cannabis, the educational attainment polygenic score did not have unique predictive power above and beyond the educational attainment phenotype.

Secondly, our sibling comparison analyses demonstrated that polygenic scores were significant predictors of SUD criteria even within families. For outcomes such as SUDs, which have considerable influences that vary among families, ruling out familial confounding is particularly important. In addition to significant sibling differences, we also found that between-family differences in *EduYears-GPS* predicted SUDs. This suggests that both the overall polygenic loading of one's family and one's relative polygenic loading within that family are important

Table 3 Sibling comparisons of substance use disorder criterion counts as a function of *EduYears-GPS* genome-wide polygenic scores.

Parameter	Alcohol use disorder criterion count (n = 3098)		Nicotine dependence criterion count (n = 2752)		Cannabis use disorder criterion count (n = 3097)	
	b	95% CI	b	95% CI	b	95% CI
Intercept	-6.88	(-8.26, -5.50)	-2.75	(-4.29, -1.21)	-1.87	(-3.00, -0.74)
Sex (female)	-0.51	(-0.57, -0.46)	-0.19	(-0.25, -0.13)	-0.40	(-0.46, -0.35)
Age (years)	2.00E-03	(-3.83E-03, 7.45E-03)	0.02	(0.01, 0.02)	-4.00E-03	(-9.90E-03, 1.45E-03)
PC1	-10.02	(-117.21, 97.16)	-186.36	(-305.61, -67.10)	-126.75	(-235.49, -18.00)
PC2	-0.50	(-53.35, 52.36)	41.30	(-17.28, 99.88)	-4.57	(-58.33, 49.18)
<i>EduYears-GPS</i> count	<b>1.78E-05</b>	<b>(1.54E-05, 2.01E-05)</b>	<b>1.37E-05</b>	<b>(1.05E-05, 1.70E-05)</b>	<b>5.47E-05</b>	<b>(4.09E-05, 6.85E-05)</b>
Cohort 2	0.34	(0.12, 0.56)	0.27	(-0.03, 0.57)	0.03	(-0.20, 0.26)
Cohort 3	0.53	(0.27, 0.79)	0.32	(-0.02, 0.66)	0.63	(0.37, 0.90)
Cohort 4	0.28	(-0.03, 0.59)	0.22	(-0.18, 0.61)	0.48	(0.17, 0.79)
<i>EduYears-GPS</i> -deviation	-19694.89 <sup>a</sup>	(-31041.68, -8348.11)	-16676.08	(-27013.63, -6338.52)	-3829.85	(-6486.43, -1173.26)
<i>EduYears-GPS</i> -mean	-13147.15	(-23328.50, -2965.80)	-29327.98	(-38688.84, -19967.11)	-2870.93	(-5437.72, -304.13)
<i>EduYears-GPS</i> -deviation $\Delta R^2$	0.17%		0.20%		0.13%	
<i>EduYears-GPS</i> -mean $\Delta R^2$	0.29%		1.89%		0.22%	

Bold type indicates estimate  $P \leq 0.05$ . <sup>a</sup>The *EduYears-GPS*-deviation effect was robust after controlling for sibling and family differences in phenotypic educational attainment (see Supporting information, Table S3). The *EduYears-GPS* thresholds for each substance were: alcohol ( $P < 0.30$ ); nicotine ( $P < 0.20$ ); cannabis ( $P < 0.01$ ). PC = principal component for genetic ancestry; cohort = dummy-coded variables indexing year of birth, defined as (1930-50); (1950-70); and (1970-2010); *EduYears-GPS* count = number of single nucleotide polymorphisms available for polygenic scoring; *EduYears-GPS*-deviation = the difference of an individual's *EduYears-GPS* from the sibling group mean; *EduYears-GPS*-mean = the sibling group mean of *EduYears-GPS*;  $\Delta R^2$  = change in  $r$ -squared.

predictors of risk for SUDs. The associations between sibling differences in polygenic scores and SUDs were attenuated somewhat after controlling for sibling differences in phenotypic educational attainment. This attenuation may reflect the relative statistical power of polygenic scores compared to the phenotypes from which they are derived, as well the likelihood that some of the effect of sibling differences in educational attainment polygenic scores is likely to be mediated through sibling differences in educational persistence, as has been documented previously [26].

These results should be considered in the context of several limitations. First, the COGA sample is enriched with individuals with SUDs, and the results may not generalize to lower-risk samples. Secondly, the sibling comparison design assumes that siblings are reared together. Not all COGA participants were asked about their living arrangements while growing up, so we could not test whether this assumption was met for all sibling groups. However, to address this concern, we restricted the analyses to the sibling groups where it was possible to determine that they grew up together, and to siblings who were born close together in time (and thus more likely to share aspects of their rearing environment compared to siblings born further apart). The pattern of effects remained significant and in the same direction in these sensitivity analyses, suggesting that the effects observed in our sibling comparisons of polygenic scores were not driven by differences in siblings' rearing environments.

Thirdly, because genetic associations can differ across ancestral groups, we focused here on the European ancestry subset of COGA because the discovery GWAS for educational attainment used a European ancestry sample. It is unknown whether the same pattern of effects would be observed in samples of non-European ancestry.

Fourthly, polygenic scores by design only capture common genetic variation. Fifthly, despite evidence for polygenic association even after controlling for family-level confounders, the polygenic scores accounted for a relatively small amount of variance. This limited predictive power cautions against incorporating polygenic scores into clinical screening or intervention efforts for substance use disorders.

As efforts to characterize how polygenic predispositions influence complex behavioral outcomes increase in popularity [16], we believe that environmentally informed designs such as sibling comparisons will become a particularly useful tool to illuminate the 'chains of risk' from genotype to phenotype. For example, sibling differences can be elaborated upon to include examination of how subtle differences in polygenic loading between siblings impact individual differences or selection into particular environments. In turn, these mediating phenotypes may be particularly actionable targets for prevention and intervention efforts.

## Declaration of interests

None.

## Acknowledgements

The Collaborative Study on the Genetics of Alcoholism (COGA), Principal Investigators (PI) B. Porjesz, V. Hesselbrock, H. Edenberg and L. Bierut, includes 11 different centers: University of Connecticut (V. Hesselbrock); Indiana University (H. J. Edenberg, J. Nurnberger Jr, T. Foroud); University of Iowa (S. Kuperman, J. Kramer); SUNY Downstate (B. Porjesz); Washington University in St Louis (L. Bierut, J. Rice, K. Bucholz, A. Agrawal); University of California at San Diego (M. Schuckit); Rutgers University (J. Tischfield, A. Brooks); Department of Biomedical and Health Informatics, The Children's Hospital of Philadelphia; Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia PA (L. Almasy), Virginia Commonwealth University (D. Dick), Icahn School of Medicine at Mount Sinai (A. Goate) and Howard University (R. Taylor). Other COGA collaborators include: L. Bauer (University of Connecticut); J. McClintick, L. Wetherill, X. Xuei, Y. Liu, D. Lai, S. O'Connor, M. Plawecki, S. Lourens (Indiana University); G. Chan (University of Iowa; University of Connecticut); J. Meyers, D. Chorlian, C. Kamarajan, A. Pandey, J. Zhang (SUNY Downstate); J.-C. Wang, M. Kapoor, S. Bertelsen (Icahn School of Medicine at Mount Sinai); A. Anokhin, V. McCutcheon, S. Saccone (Washington University); J. Salvatore, E. Aliev, B. Cho (Virginia Commonwealth University); and Mark Kos (University of Texas Rio Grande Valley). A. Parsian and M. Reilly are the NIAAA Staff Collaborators. We continue to be inspired by our memories of Henri Begleiter and Theodore Reich, founding PI and Co-PI of COGA, and also owe a debt of gratitude to other past organizers of COGA, including Ting-Kai Li, P. Michael Conneally, Raymond Crowe and Wendy Reich for their critical contributions. This national collaborative study is supported by NIH Grant U10AA008401 from the National Institute on Alcohol Abuse and Alcoholism (NIAAA) and the National Institute on Drug Abuse (NIDA). Additional support for this project comes from K01AA024152 (J.E. S.); K02DA032573 (A.A.); and R01DA040411 (E.C.J.). Funding support for GWAS genotyping performed at the Johns Hopkins University Center for Inherited Disease Research was provided by the National Institute on Alcohol Abuse and Alcoholism, the NIH GEI (U01HG004438), and the NIH contract 'High throughput genotyping for studying the genetic contributions to human disease' (HHSN268200782096C). GWAS genotyping was also performed at the Genome Technology Access Center in the Department of Genetics at Washington University School of Medicine, which is partially supported by NCI Cancer Center Support Grant no. P30 CA91842 to the



Siteman Cancer Center and by ICTS/CTSA Grant# UL1R024992 from the National Center for Research Resources (NCRR), a component of the National Institutes of Health (NIH), and NIH Roadmap for Medical Research.

### Author's affiliations

Department of Psychology, Virginia Commonwealth University, Richmond, VA, USA,<sup>1</sup> Virginia Institute for Psychiatric and Behavioral Genetics, Virginia Commonwealth University, Richmond, VA, USA,<sup>2</sup> Department of Business Administration, Karabuk University, Karabuk, Turkey,<sup>3</sup> Department of Psychology, Arizona State University, Tempe, AZ, USA,<sup>4</sup> Department of Psychiatry, Washington University, St Louis, MO, USA,<sup>5</sup> Department of Genetics, University of Pennsylvania, Philadelphia, PA, USA,<sup>6</sup> Department of Biomedical and Health Informatics, Children's Hospital of Philadelphia, Philadelphia, PA, USA,<sup>7</sup> Department of Psychiatry, University of Connecticut School of Medicine, Farmington, CT, USA,<sup>8</sup> Department of Biochemistry and Molecular Biology, Indiana University, Indianapolis, IN, USA,<sup>9</sup> Department of Psychiatry, SUNY Downstate Medical Center, Brooklyn, NY, USA,<sup>10</sup> Department of Psychiatry, University of California—San Diego, La Jolla, CA, USA,<sup>11</sup> Department of Genetics and the Human Genetics Institute of New Jersey, Piscataway, NJ, USA,<sup>12</sup> Department of Medical and Molecular Genetics, Indiana University, Indianapolis, IN, USA,<sup>13</sup> Department of Human and Molecular Genetics, Virginia Commonwealth University, Richmond, VA, USA<sup>14</sup> and College Behavioral and Emotional Health Institute, Virginia Commonwealth University, Richmond, VA, USA<sup>15</sup>

### References

- Esch P, Bocquet V, Pull C., Couffignal S., Lehnert T., Graas M., *et al.* The downward spiral of mental disorders and educational attainment: a systematic review on early school leaving. *BMC Psychiatry* 2014; **14**: 237.
- Townsend L., Flisher A. J., King G. A systematic review of the relationship between high school dropout and substance use. *Clin Child Fam Psychol Rev* 2007; **10**: 295–317.
- Kessler R. C., Foster C. L., Saunders W. B., Stang P. E. Social consequences of psychiatric disorders, I: Educational attainment. *Am J Psychiatry* 1995; **152**: 1026–32.
- Stinson F. S., Ruan W. J., Pickering R., Grant B. F. Cannabis use disorders in the USA: prevalence, correlates and co-morbidity. *Psychol Med* 2006; **36**: 1447–60.
- Breslau J., Lane M., Sampson N., Kessler R. C. Mental disorders and subsequent educational attainment in a US national sample. *J Psychiatr Res* 2008; **42**: 708–16.
- Verweij K. J., Huizink A. C., Agrawal A., Martin N. G., Lynskey M. T. Is the relationship between early-onset cannabis use and educational attainment causal or due to common liability? *Drug Alcohol Depend* 2013; **133**: 580–6.
- Breslau J., Miller E., Chung W. J., Schweitzer J. B. Childhood and adolescent onset psychiatric disorders, substance use, and failure to graduate high school on time. *J Psychiatr Res* 2011; **45**: 295–301.
- Martin M. J., Conger R. D., Sitnick S. L., Masarik A. S., Forbes E. E., Shaw D. S. Reducing risk for substance use by economically disadvantaged young men: positive family environments and pathways to educational attainment. *Child Dev* 2015; **86**: 1719–37.
- Fothergill K. E., Ensminger M. E., Green K. M., Crum R. M., Robertson J., Juon H. S. The impact of early school behavior and educational achievement on adult drug use disorders: a prospective study. *Drug Alcohol Depend* 2008; **91**: 191–9.
- Green K. M., Zembrak K. A., Fothergill K. E., Robertson J. A., Ensminger M. E. Childhood and adolescent risk factors for comorbid depression and substance use disorders in adulthood. *Addict Behav* 2012; **37**: 1240–7.
- Latvala A., Dick D. M., Tuulio-Henriksson A., Suvisaari J., Viken R. J., Rose R. J., *et al.* Genetic correlation and gene–environment interaction between alcohol problems and educational level in young adulthood. *J Stud Alcohol Drugs* 2011; **72**: 210–20.
- Bulik-Sullivan B., Finucane H. K., Anttila V., Gusev A., Day F. R., Loh P. R., *et al.* An atlas of genetic correlations across human diseases and traits. *Nat Genet* 2015; **47**: 1236–41.
- Walters R. K., Polimanti R., Johnson E. C., Mcclintick J. N., Adams M. J., Adkins A. E., *et al.* Trans-ancestral GWAS of alcohol dependence reveals common genetic underpinnings with psychiatric disorders. *Nat Neurosci* 2018; **21**: 1656–69.
- Bergen S. E., Gardner C. O., Aggen S. H., Kendler K. S. Socio-economic status and social support following illicit drug use: causal pathways or common liability? *Twin Res Hum Genet* 2008; **11**: 266–74.
- Grant J. D., Scherrer J. F., Lynskey M. T., Agrawal A., Duncan A. E., Haber J. R., *et al.* Associations of alcohol, nicotine, cannabis, and drug use/dependence with educational attainment: evidence from cotwin-control analyses. *Alcohol Clin Exp Res* 2012; **36**: 1412–20.
- Martin A. R., Daly M. J., Robinson E. B., Hyman S. E., Neale B. M. Predicting polygenic risk of psychiatric disorders. *Biol Psychiatry* 2019; **86**: 97–109.
- Maier R. M., Visscher P. M., Robinson M. R., Wray N. R. Embracing polygenicity: a review of methods and tools for psychiatric genetics research. *Psychol Med* 2017; **48**: 1055–67.
- Bogdan R., Baranger D. A. A., Agrawal A. Polygenic risk scores in clinical psychology: bridging genomic risk to individual differences. *Annu Rev Clin Psychol* 2018; **14**: 119–57.
- Dudbridge F. Polygenic epidemiology. *Genet Epidemiol* 2016; **40**: 268–72.
- Lee J. J., Wedow R., Okbay A., Kong E., Maghziyan O., Zacher M., *et al.* Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nat Genet* 2018; **50**: 1112–21.
- D'Onofrio B. M., Lahey B. B., Turkheimer E., Lichtenstein P. Critical need for family-based, quasi-experimental designs in integrating genetic and social science research. *Am J Public Health* 2013; **103**: 46–55.
- Lahey B. B., D'Onofrio B. M. All in the family: comparing siblings to test causal hypotheses regarding environmental influences on behavior. *Curr Dir Psychol Sci* 2010; **19**: 319–23.
- Rutter M. Proceeding from observed correlation to causal inference: the use of natural experiments. *Perspect Psychol Sci* 2007; **2**: 377–95.
- Donovan S. J., Susser E. Commentary: advent of sibling designs. *Int J Epidemiol* 2011; **40**: 345–9.
- Belsky D. W., Moffitt T. E., Corcoran D. L., Domingue B., Harrington H., Hogan S., *et al.* The genetics of success: how single-nucleotide polymorphisms associated with educational attainment relate to life-course development. *Psychol Sci* 2016; **27**: 957–72.
- Domingue B. W., Belsky D., Conley D., Harris K. M., Boardman J. D. Polygenic influence on educational attainment: new evidence from the National Longitudinal Study of adolescent to adult health. *AERA Open* 2015; **1**: 1–13.
- Selzam S., Krapohl E., Von Stumm S., O'Reilly P. F., Rimfeld K., Kovas Y., *et al.* Predicting educational achievement from DNA. *Mol Psychiatry* 2017; **22**: 267–72.

28. Galea S., Nandi A., Vlahov D. The social epidemiology of substance use. *Epidemiol Rev* 2004; **26**: 36–52.
29. Begleiter H., Reich T., Hesselbrock V., Porjesz B., Li T. K., Schuckit M. A., et al. The collaborative Study on the genetics of alcoholism. *Alcohol Health Res W* 1995; **19**: 228–36.
30. Bucholz K. K., McCutcheon V. V., Agrawal A., Dick D. M., Hesselbrock V. M., Kramer J. R., et al. Comparison of parent, peer, psychiatric, and cannabis use influences across stages of offspring alcohol involvement: evidence from the COGA prospective study. *Alcohol Clin Exp Res* 2017; **41**: 359–68.
31. Lai D., Wetherill L., Bertelsen S., Carey C. E., Kamarajan C., Kapoor M., et al. Genome-wide association studies of alcohol dependence, DSM-IV criterion count and individual criteria. *Genes Brain Behav* 2019; **18**: e12579.
32. Cardon L. R., Palmer L. J. Population stratification and spurious allelic association. *Lancet* 2003; **361**: 598–604.
33. Zheng X., Levine D., Shen J., Gogarten S. M., Laurie C., Weir B. S. A high-performance computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* 2012; **28**: 3326–8.
34. Grant B. F., Goldstein R. B., Saha T. D., Chou S. P., Jung J., Zhang H., et al. Epidemiology of DSM-5 alcohol use disorder: results from the National Epidemiologic Survey on alcohol and related conditions III. *JAMA Psychiatry* 2015; **72**: 757–66.
35. Hasin D. S., Kerridge B. T., Saha T. D., Huang B., Pickering R., Smith S. M., et al. Prevalence and correlates of DSM-5 cannabis use disorder, 2012–2013: findings from the National Epidemiologic Survey on alcohol and related conditions—III. *Am J Psychiatry* 2016; **173**: 588–99.
36. Breslau N., Johnson E. O., Hiripi E., Kessler R. Nicotine dependence in the United States: prevalence, trends and smoking persistence. *Arch Gen Psychiatry* 2001; **58**: 810–6.
37. US Census Bureau. Current Population Survey, 2018 Annual Social and Economic Supplement. Suitland, MD; 2018.
38. Hesselbrock M., Easton C., Bucholz K. K., Schuckit M., Hesselbrock V. A validity study of the SSAGA—a comparison with the SCAN. *Addiction* 1999; **94**: 1361–70.
39. Bucholz K. K., Cadoret R., Cloninger C. R., Dinwiddie S. H., Hesselbrock V. M., Nurnberger J. L., et al. A new, semi-structured psychiatric interview for use in genetic linkage studies: a report on the reliability of the SSAGA. *J Stud Alcohol* 1994; **55**: 149–58.
40. American Psychiatric Association *Diagnostic and Statistical Manual of Mental Disorders, 5th edn*. Arlington, VA: American Psychiatric Publishing; 2013.
41. Heatherton T. F., Kozłowski L. T., Frecker R. C., Fagerstrom K. O. The Fagerstrom test for nicotine dependence: a revision of the Fagerstrom tolerance questionnaire. *Br J Addict* 1991; **86**: 1119–27.
42. Purcell S., Neale B., Todd-Brown K., Thomas L., Ferreira M., Bender D., et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* 2007; **81**: 559–75.
43. Pinheiro J., EISPACK, Heisterkamp S., Van Willigen B., R Core Team. Linear and nonlinear mixed effects models. 2018. <https://CRAN.R-project.org/package=nlme>
44. R Development Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria; 2014.
45. Euesden J., Lewis C. M., O'Reilly P. E. PRSice: polygenic risk score software. *Bioinformatics* 2015; **31**: 1466–8.
46. Barton K. MuMIn: Multi-Model Inference. R package version 1.15.6; 2016. <https://cran.r-project.org/web/packages/MuMIn/index.html>
47. Krapohl E., Rimfeld K., Shakeshaft N. G., Trzaskowski M., McCmillan A., Pingault J. B., et al. The high heritability of educational achievement reflects many genetically influenced traits, not just intelligence. *Proc Natl Acad Sci USA* 2014; **111**: 15273–8.
48. Okbay A., Beauchamp J. P., Fontana M. A., Lee J. J., Pers T. H., Rietveld C. A., et al. Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* 2016; **533**: 539–42.
49. Hagenaars S. P., Harris S. E., Davies G., Hill W. D., Liewald D. C., Ritchie S. J., et al. Shared genetic aetiology between cognitive functions and physical and mental health in UK biobank (N=112 151) and 24 GWAS consortia. *Mol Psychiatry* 2016; **21**: 1624–32.
50. De Zeeuw E. L., Van Beijsterveldt C. E., Glasner T. J., Bartels M., Ehli E. A., Davies G. E., et al. Polygenic scores associated with educational attainment in adults predict educational achievement and ADHD symptoms in children. *Am J Med Genet B Neuropsychiatr Genet* 2014; **165B**: 510–20.

### Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

#### Table S1

Associations between *EduYears-GPS* genome-wide polygenic scores across multiple thresholds and substance use disorder criterion counts (presented as t-values)

**Table S2** Associations between *EduYears-GPS* genome-wide polygenic scores and substance use disorder criterion counts, controlling for educational attainment

**Table S3** Sibling comparisons of substance use disorder criterion counts as a function of *EduYears-GPS* genome-wide polygenic scores, controlling for sibling and family differences in educational attainment.

**Table S4** Sibling comparisons of substance use disorder criterion counts as a function of *EduYears-GPS* genome-wide polygenic scores for siblings known to have same childhood/adolescent household structure.

**Table S5** Sibling comparisons of substance use disorder criterion counts as a function of *EduYears-GPS* genome-wide polygenic scores for siblings born within 3 years of one another.

**Table S6** Sibling comparisons of substance use disorder criterion counts as a function of *EduYears-GPS* genome-wide polygenic scores, excluding sibling groups that included monozygotic twins.

**Figure S1** Illustration of sibling comparison sample sizes as a function of filtering criteria employed for the primary sample and the sensitivity analyses.