

# Ultra-rare disruptive and damaging mutations influence educational attainment in the general population

Andrea Ganna<sup>1-4,15</sup>, Giulio Genovese<sup>2,3,5,15</sup>, Daniel P Howrigan<sup>1-3</sup>, Andrea Byrnes<sup>1-3</sup>, Mitja I Kurki<sup>1-3,6</sup>, Seyedeh M Zekavat<sup>2,7</sup>, Christopher W Whelan<sup>2,3,5</sup>, Mart Kals<sup>8,9</sup>, Michel G Nivard<sup>10</sup>, Alex Bloemendal<sup>1-3</sup>, Jonathan M Bloom<sup>1-3</sup>, Jacqueline I Goldstein<sup>1-3</sup>, Timothy Poterba<sup>1-3</sup>, Cotton Seed<sup>1-3</sup>, Robert E Handsaker<sup>2,3,5</sup>, Pradeep Natarajan<sup>2,7</sup>, Reedik Mägi<sup>8</sup>, Diane Gage<sup>3</sup>, Elise B Robinson<sup>1-3</sup>, Andres Metspalu<sup>8</sup>, Veikko Salomaa<sup>11</sup>, Jaana Suvisaari<sup>11</sup>, Shaun M Purcell<sup>12,13</sup>, Pamela Sklar<sup>12,13</sup>, Sekar Kathiresan<sup>2,7</sup>, Mark J Daly<sup>1-3</sup>, Steven A McCarroll<sup>2,3,5</sup>, Patrick F Sullivan<sup>4,14</sup>, Aarno Palotie<sup>1,3,6</sup>, Tõnu Esko<sup>2,8</sup>, Christina M Hultman<sup>4</sup> & Benjamin M Neale<sup>1-3</sup>

**Disruptive, damaging ultra-rare variants in highly constrained genes are enriched in individuals with neurodevelopmental disorders. In the general population, this class of variants was associated with a decrease in years of education (YOE). This effect was stronger among highly brain-expressed genes and explained more YOE variance than pathogenic copy number variation but less than common variants. Disruptive, damaging ultra-rare variants in highly constrained genes influence the determinants of YOE in the general population.**

Educational attainment, measured by the highest number of YOE attained, is a complex trait influenced by public policy<sup>1</sup>, economic resources<sup>2</sup> and many heritable traits, including cognitive abilities and behavior<sup>3</sup>. Importantly, YOE is positively associated with healthy behaviors and lower rates of chronic diseases<sup>4</sup>. Genome-wide association study (GWAS) meta-analyses have identified 162 genome-wide significant loci for YOE<sup>5</sup>. The additive heritability of YOE explained by common genetics variants has been estimated at 21% (95% confidence intervals (CI): 11–31%)<sup>6</sup>, which is approximately half of the total heritability estimated from twin studies (40%; 95% CI: 35–44%)<sup>7</sup>. Rare to ultra-rare exonic variants might account for some of the heritability currently not captured by GWAS<sup>8</sup>.

Recent studies of intellectual disability, autism and schizophrenia have shed light on the impact of *de novo* and ultra-rare variants (URVs: variants that are observed only once (singletons) in the study and not observed in

60,706 exomes sequenced in the Exome Aggregation Consortium<sup>9</sup>) on the genetic architecture of these disorders<sup>10–12</sup>, showing a specific enrichment in highly constrained (HC: genes intolerant to loss-of-function or missense mutations, i.e., having a probability of being loss-of-function intolerant ( $P_{LI}$ ) > 0.9). Moreover, emerging evidence suggests that *de novo* loss-of-function mutations are associated with reduced adaptive functioning in individuals who have not been diagnosed with autism<sup>13</sup>.

We tested the hypothesis that a burden of URVs in HC genes is associated with YOE in 14,133 individuals participating in four studies from three Northern European countries: Sweden, Estonia and Finland. Of these, 5,047 individuals have been diagnosed with schizophrenia. The average numbers of YOE were 13.1, 13.6 and 11.8 in Swedish, Estonian and Finnish participants, respectively. These differences are partially explained by different age and sex distributions, as well as by different methods used to measure educational attainment (**Supplementary Table 1**). We observed lower YOE among men compared to women (12.8 versus 13.2 years,  $P = 4.8 \times 10^{-12}$ ) and older individuals compared to younger (0.8 month less of education for each additional year of age,  $P < 1 \times 10^{-15}$ ; **Supplementary Table 2**).

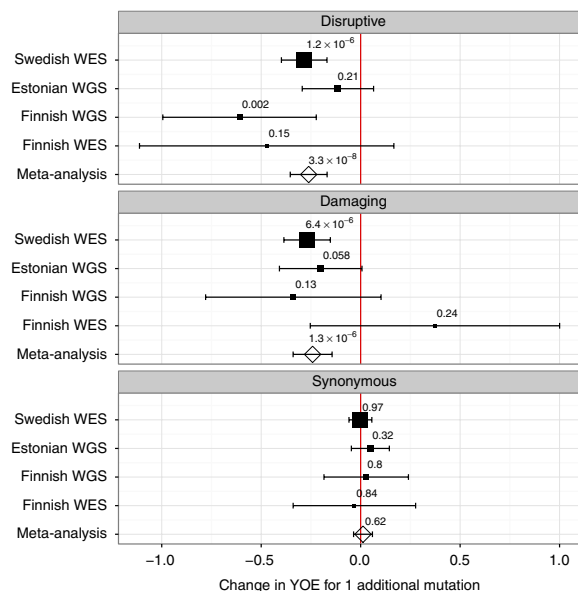
We developed a new software package, Hail, to efficiently perform quality control, annotation and analysis of large-scale sequencing data (Online Methods). We identified URVs in HC genes using whole-exome sequencing (WES) data ( $n = 11,431$  individuals) and protein-coding regions in high-coverage whole-genome sequencing (WGS) data ( $n = 2,702$  individuals). The primary reason for focusing on URVs in HC genes was to maximize the expected deleteriousness of the variants included (due to purifying selection).

Within the set of URVs in HC genes we defined variants that were (i) disruptive: putative loss-of-function variants including premature stop codons, essential splice site mutations and frame-shift indels; (ii) damaging: missense variants classified as damaging by seven different *in silico* prediction algorithms (Online Methods) and (iii) negative control: synonymous variants not predicted to change the encoded protein. We observed one or more of such mutations in 25%, 24% and 78% of individuals, respectively (**Supplementary Table 3**). Principal components of genetic data showed that individuals within each study were of similar ancestry (**Supplementary Fig. 1**).

On average (**Fig. 1**), we observed a 3.1-month reduction in YOE for each disruptive mutation (95% CI:  $-4.3, -2.0$ ;  $P = 3.3 \times 10^{-8}$ ) and a similar effect for damaging mutations (2.9 months fewer YOE; 95% CI:  $-4.1, -1.7$ ;  $P = 1.3 \times 10^{-6}$ ). Furthermore, each additional disruptive mutation on average reduced the chance of going to college by 14% (odds ratio = 0.86; 95% CI: 0.78, 0.95;  $P = 0.0017$ ). These results were consistent

<sup>1</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts, USA. <sup>2</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA. <sup>3</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA. <sup>4</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden. <sup>5</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts, USA. <sup>6</sup>Institute for Molecular Medicine Finland, FIMM, University of Helsinki, Helsinki, Finland. <sup>7</sup>Center for Human Genetic Research and Cardiovascular Research Center, Massachusetts General Hospital, Boston, Massachusetts, USA. <sup>8</sup>Estonian Genome Center, University of Tartu, Tartu, Estonia. <sup>9</sup>Institute of Mathematics and Statistics, University of Tartu, Tartu, Estonia. <sup>10</sup>Department of Biological Psychology, VU University Amsterdam, Amsterdam, the Netherlands. <sup>11</sup>Department of Health, THL-National Institute for Health and Welfare, Helsinki, Finland. <sup>12</sup>Division of Psychiatric Genomics, Department of Psychiatry, Icahn School of Medicine at Mount Sinai, New York, New York, USA. <sup>13</sup>Institute for Genomics and Multiscale Biology, Icahn School of Medicine at Mount Sinai, New York, New York, USA. <sup>14</sup>Departments of Genetics and Psychiatry, University of North Carolina, Chapel Hill, North Carolina, USA. <sup>15</sup>These authors contributed equally to this work. Correspondence should be addressed to A.G. ([aganna@broadinstitute.org](mailto:aganna@broadinstitute.org)).

Received 12 June; accepted 7 September; published online 3 October 2016; doi:10.1038/nn.4404



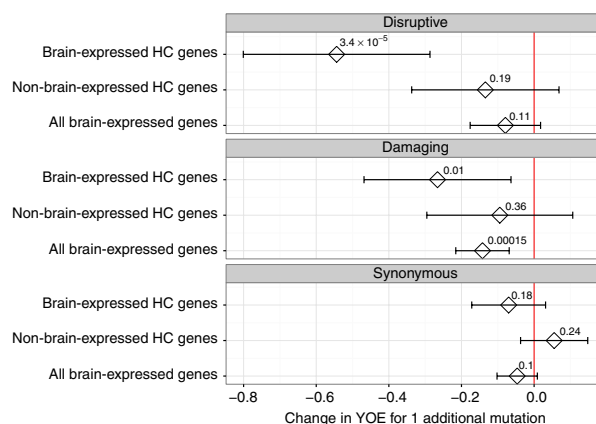
**Figure 1** Association between number of disruptive, damaging and synonymous URVs in HC genes and YOE. Disruptive and damaging URVs but not synonymous URVs are significantly associated with reduced YOE. The size of the squares is proportional to the size of the study. Bars are 95% confidence intervals. All estimates were obtained from a linear regression model. Meta-analysis results were obtained using a fixed-effect approach.

when using a mixed linear model approach to correct for population stratification in the Finnish and Estonian samples with WGS data (2.4 months fewer YOE; 95% CI:  $-4.3, -0.95$ ;  $P = 0.014$ ,  $n = 2,702$ ).

The negative association between URVs and YOE remained consistent when we examined the control cohort and schizophrenia cohort separately (**Supplementary Fig. 2**). Furthermore, the effect remained consistent when excluding individuals diagnosed with a neurodevelopmental disorder (i.e., schizophrenia, bipolar disorder, autism, intellectual disability and Asperger's syndrome), as identified via linkage with the Swedish national inpatient registry (**Supplementary Fig. 3**). We did not observe any significant association when we restricted our analysis to synonymous variants in HC genes ( $P = 0.62$ ) or disruptive mutations in unconstrained genes ( $P = 0.73$ ).

We used gene-expression data to determine whether restricting our analysis to genes enriched for brain expression concentrated our URVs burden signal. Specifically, we used the Genotype-Tissue Expression consortium data<sup>14</sup> to identify the top 20% brain-expressed HC genes. The intersection between HC and brain-expressed genes ( $n = 683$  genes and 313 genes for disruptive and damaging URVs, respectively) more than doubled the impact on YOE (6.5 months fewer YOE per each additional disruptive variant; 95% CI:  $-9.6, -3.4$ ;  $P = 3.4 \times 10^{-5}$ ; **Fig. 2**). When using increasingly liberal thresholds for defining genes enriched for brain-expression, we saw a consistent decrease in the association (**Supplementary Fig. 4**). The association was not significant when considering disruptive URVs in non-brain-enriched HC genes or all brain-enriched genes ( $P = 0.19$  and  $P = 0.11$ , respectively). We further examined a subset of genes for which basal gene-expression was at least two-fold higher in the brain compared to other tissues (brain-enriched HC genes; **Supplementary Fig. 5**). Although the impact on YOE was higher for brain-enriched HC genes than for non-brain-enriched HC genes, the signal was specific to disruptive variants. Overall, this approach was less effective in identifying a HC gene subset impacting YOE.

To place disruptive and damaging URVs into context, we also examined the impact of previously reported genetic influences on YOE,



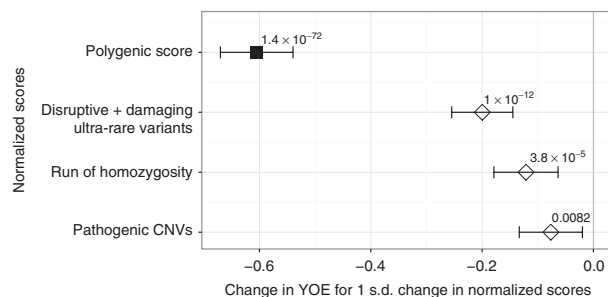
**Figure 2** Association between numbers of disruptive, damaging and synonymous URVs for different gene sets. The intersection between HC and brain-expressed genes yields the strongest reduction in YOE. We only report the meta-analysis results ( $n = 14,133$ ). Bars are 95% confidence intervals. All estimates were obtained from a linear regression model and combined using fixed-effect meta-analysis.

including a polygenic score from common variants<sup>5</sup>, runs of homozygosity<sup>15</sup> and a burden of rare pathogenic copy number variants (CNVs)<sup>16</sup>. We sought to establish whether or not these different forms of genetic variation act independently on YOE. For this purpose we defined four scores: (i) a polygenic score including all the independent single nucleotide polymorphisms with  $P < 1$  for association with YOE (as this threshold has been shown to maximize variance explained in YOE) in a large GWAS consortia of YOEs<sup>5</sup>, (ii) the summed length of all runs of homozygosity, (iii) the burden of disruptive and damaging URVs in HC genes and (iv) the burden of a self-curated list of pathogenic CNVs from the literature (**Supplementary Table 4**). The polygenic score was only calculated for the Swedish samples ( $n = 10,644$ ), since the other three studies were included in the original GWAS of YOE.

We first explored the association between each genetic score and YOE separately. The strongest change in YOE was observed among CNV carriers (7.6 months fewer YOE; 95% CI:  $-13.7, -1.5$ ;  $P = 0.015$ ). However, these events were rare in the population (161 carriers among 11,999 individuals with CNV measured). We then fit the four normalized scores in the same regression model to assess the relative contribution of each genetic class to YOE. All four scores were independently associated with YOE (**Fig. 3**). The polygenic score showed the strongest association in standard deviation from the mean, explaining the largest proportion of the variability in YOE (2.9% versus 0.4% for the ultra-rare variants, 0.2% for runs of homozygosity and 0.1% for pathogenic CNVs). We further evaluated whether the association between the polygenic score and YOE changed in individuals with and without disruptive or damaging URVs or CNVs. We found that the polygenic score was more strongly associated with YOE in individuals without disruptive or damaging URVs or CNVs (8.2 versus 6.2 more months of YOE for 1 standard deviation increase in the polygenic score;  $P(\text{interaction}) = 0.007$ ; **Supplementary Fig. 6**).

We sought to identify individual genes driving the observed association between disruptive and damaging URVs and YOE. Using a gene-based burden test implemented in SKAT (ref. 17) and an exome-wide significance threshold of  $1 \times 10^{-6}$ , we identified no statistically significantly associated genes (**Supplementary Fig. 7**, upper panels). Similar results were observed when we included all variants with minor allele frequency  $< 0.05\%$ , rather than only URVs (**Supplementary Fig. 7**, lower panels).

In this study we focused on YOE, a phenotype that is relatively easy to collect in large samples and that has a strong genetic



**Figure 3** Association between each of the normalized scores (polygenic, runs of homozygosity, URVs and pathogenic CNVs) and YOE. The results presented are from meta-analysis of Swedish WES, Estonian WGS and Finnish WGS studies ( $n = 13,353$ ), except for the polygenic score, which is calculated only in the Swedish WES study ( $n = 10,651$ ). Note that we plotted 1 polygenic score to obtain a negative association with YOE. The horizontal bars represent 95% confidence intervals. All the estimates were obtained from a linear regression model and combined using fixed-effect meta-analysis.

correlation with intelligence and cognitive function<sup>6,18</sup>. We integrated WGS, WES and array data on more than 14,000 individuals and described the impact of URVs disrupting HC genes on YOE. This class of variants has been previously associated with autism<sup>3</sup> and schizophrenia<sup>4</sup>, but the impact on YOE in the general population has not, to our knowledge, been described before. Here we show that disruptive and damaging URVs in HC genes are likely to affect factors underlying education attainment among individuals not diagnosed with psychiatric or neurodevelopmental disorders. Exploring the extent to which this association is mediated by cognitive-related determinants of YOE or by other noncognitive factors will require studies integrating detailed cognitive, psychological and personality measurements.

Similarly to genetic analyses of schizophrenia<sup>10</sup> and autism, the majority of the signal lies in genes highly expressed in brain. This observation does not exclude the existence of causal mutations outside this gene class, but it suggests that strong-acting mutations are heavily concentrated within these genes. Furthermore, we showed that disruptive and damaging URVs in HC genes, common variants associated with YOE, runs of homozygosity and pathogenic CNVs all act on cognitive function or personality traits ultimately reflected in the educational attainment of our study participants. This effect was not simply additive. We identified a modest but significant interaction between the polygenic score and the presence of URVs or CNVs ( $P = 0.007$ ). Whether this observation is driven by the interplay of partially overlapping pathways between common and rare variants or by genotype–phenotype heterogeneity (for example, common and rare variants impacting different subsets of individuals) will be a matter for future investigation.

We found that, on average, each additional disruptive URVs in HC genes results in a 3.1-month reduction in YOE. This effect is likely to be a mixture of variants with larger effect and variants that are not associated with YOE. The polygenic score based on common-variant effect sizes estimated from a much larger cohort of 405,072 individuals explained a larger fraction of the YOE. This is not surprising, given that common variants are expected to have the largest contribution to heritable variation in most complex traits<sup>19</sup>.

The prioritization approaches used to select variants contributing to the score from common variants and to select those contributing to the score from rare variants are different. The former uses estimates of the association with YOE, and the proportion of variance explained by the score is likely to improve as the sample size used to originate these estimates increases. The latter uses *in silico* prediction of the variants' functional effect coupled with population genetics expectations built

on the mutation rate. As with the common-variant score, we expect that the score based on URVs in selected gene sets will continue to improve in predictive validity of YOE as a more precise characterization of which genes and genomic regions are associated with YOE emerges.

Our study could not detect disruptive or damaging mutations in a given gene as being unequivocally associated with YOE; however, as sample sizes increase, specific genes will emerge. Nevertheless, our proof-of-concept work shows that a wide range of genetic variation, from URVs and CNVs to common variants, influences determinants of YOE in the population.

## METHODS

Methods, including statements of data availability and any associated accession codes and references, are available in the [online version of the paper](#).

**Accession codes.** dbGAP: [pht000473.v2.p2](#).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

We thank R. Walters for discussions. A.G. is supported by the Knut and Alice Wallenberg Foundation (2015.0327) and the Swedish Research Council (2016-00250). M.G.N. is supported by the Royal Netherlands Academy of Science Professor Award (PAH/6635) to Dorret I. Boomsma. V.S. was supported by the Finnish Foundation for Cardiovascular Research. This study was supported by grants from the National Human Genome Research Institute (U54 HG003067 and R01 HG006855); the National Institute of Mental Health (1U01MH105666-01 and 1R01MH101244-02); the National Institute of Diabetes and Digestive and Kidney Disease (1U54DK105566-02); the Stanley Center for Psychiatric Research; the Alexander and Margaret Stewart Trust; the National Institutes of Mental Health (R01 MH077139 and RC2 MH089905); the Sylvan C. Herman Foundation; EU H2020 grants 692145, 676550 and 654248; Estonian Research Council Grant IUT20-60, NIASC, EIT-Health; NIH-BMI Grant No. 2R01DK075787-06A1; and by the EU through the European Regional Development Fund (Project No. 2014-2020.4.01.15-0012 GENTRANSMED).

## AUTHOR CONTRIBUTIONS

A.G. and G.G. designed the study, performed the analysis and wrote the manuscript. B.M.N. supervised the project. D.H., A. Byrnes, M.I.K., S.M.Z., C.W.W., M.K., M.G.N., P.N. and R.M. performed the analyses. A. Bloemendal, J.M.B., J.I.G., T.P., C.S. and R.E.H. developed and provided computational tools. D.G. provide data management support. E.B.R., A.M., V.S., J.S., S.M.P., P.S., S.K., M.J.D., S.A.M., P.F.S., A.P., T.E. and C.M.H. collected and provided the data.

## COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

- McLendon, M.K. & Perna, L.W. *Ann. Am. Acad. Pol. Soc. Sci.* **655**, 6–15 (2014).
- Haveman, R. & Wolfe, B. *J. Econ. Lit.* **33**, 1829–1878 (1995).
- Krapohl, E. *et al. Proc. Natl. Acad. Sci. USA* **111**, 15273–15278 (2014).
- Cutler, D.M. & Lleras-Muney, A. Education and health: evaluating theories and evidence. *National Bureau of Economic Research Working Paper Series* 12352 (2006).
- Okbay, A. *et al. Nature* **533**, 539–542 (2016).
- Marioni, R.E. *et al. Intelligence* **44**, 26–32 (2014).
- Branigan, A.R., McCallum, K.J. & Freese, J. *Soc. Forces* **92**, 109–140 (2013).
- Zuk, O. *et al. Proc. Natl. Acad. Sci. USA* **111**, E455–E464 (2014).
- Lek, M. *et al. Nature* **536**, 285–291 (2016).
- Genovese, G. *et al. Nat. Neurosci.* <http://dx.doi.org/10.1038/nn.4402> (2016).
- Gillissen, C. *et al. Nature* **511**, 344–347 (2014).
- Iossifov, I. *et al. Nature* **515**, 216–221 (2014).
- Robinson, E.B. *et al. Nat. Genet.* **48**, 552–555 (2016).
- GTEx Consortium. *Nat. Genet.* **45**, 580–585 (2013).
- Joshi, P.K. *et al. Nature* **523**, 459–462 (2015).
- Stefansson, H. *et al. Nature* **505**, 361–366 (2014).
- Wu, M.C. *et al. Am. J. Hum. Genet.* **89**, 82–93 (2011).
- Davies, K. *Br. J. Nurs.* **15**, 252–256 (2006).
- Yang, J. *et al. Nat. Genet.* **47**, 1114–1120 (2015).



## ONLINE METHODS

**Study description and selection.** In this study we used epidemiological studies with YOE and exome or whole-genome sequencing information available; no formal power calculation was done.

Ethical committees in Sweden, Estonia and Finland approved all procedures and all subjects provided written informed consent (or legal guardian consent and subject assent).

**Sweden WES.** A total of 12,384 blood-derived DNA samples from Swedish research participants were collected from 2005 to 2013. Psychiatric cases with a diagnosis of schizophrenia were ascertained from the Swedish National Hospital Discharge Register. The register is complete from 1987 and augmented by psychiatric data from 1973–1986. It contains dates and ICD discharge diagnoses (World Health Organization, 1992) for each hospitalization and captures the clinical diagnosis made by the attending physician. Case inclusion criteria:  $\geq 2$  hospitalizations with discharge diagnoses of schizophrenia, both parents born in Scandinavia and age  $\geq 18$  years. Case exclusion criteria: hospital register diagnosis of any medical or psychiatric disorder mitigating a confident diagnosis of schizophrenia as determined by expert review; this removed 3.4% of eligible cases due to the primacy of another psychiatric disorder (0.9%), a general medical condition (0.3%) or uncertainties in the Hospital Discharge Register (for example, contiguous admissions with brief total duration, 2.2%). The validity of this case definition of schizophrenia is strongly supported as described in ref. 20. Controls were selected at random from Swedish population registers. Control inclusion criteria: never hospitalized for schizophrenia or bipolar disorder (given evidence of genetic overlap with schizophrenia), both parents born in Scandinavia and age  $\geq 18$  years.

**Estonia WGS.** Estonian WGS samples are a subset of the Estonian Biobank of the Estonian Genome Center at the University of Tartu<sup>21</sup>. It is a population-based biobank, containing almost 52,000 samples of the adult population (age  $\geq 18$  years), which closely reflects the age, sex and geographical distribution of the Estonian population. All subjects have been recruited randomly by general practitioners or physicians in hospitals throughout the country. The participants donated blood samples for DNA, white blood cells and plasma tests and filled out the computer-assisted personal interview (CAPI).

In total, 2,300 geographically diverse samples have whole genome sequencing data, selected randomly by county of birth.

**Finnish WES and Finnish WGS.** All of the Finnish individuals are part of the FINRISK cohort, a national survey on risk factors of chronic and noncommunicable diseases in Finland<sup>22</sup>. The survey has been conducted every five years since 1972 in randomly selected, representative population samples from different parts of Finland. All of the samples are from FINRISK 1992, 1997, 2002 and 2007 surveys.

Finnish WES mainly includes individuals that are part of an IBD case-control study, in which controls were selected to have a high IBD polygenic risk score<sup>23</sup>.

Finnish WGS includes schizophrenia cases and controls selected using nationwide hospital discharge registry and/or nationwide medicine reimbursement registry, in which all psychosis cases or psychosis medication purchases are systematically recorded. Controls were selected to have high polygenic risk score for schizophrenia<sup>24</sup>.

**Phenotype definition.** We matched the original educational categories with the International Standard Classification of Education (ISCED), as described in **Supplementary Table 1**. Thereafter we used the equivalent of United States years of schooling to obtain the YOE. Going to college was defined as having an ISCED category  $> 4$ .

To remove potential bias introduced by uncompleted education, we excluded all the individuals younger than 30 years at the time of sample collection. For the Estonian and Finnish samples, we used self-report data; whereas for the Swedish sample, we obtained YOE from the national registries. YOE was approximately normally distributed.

**Sequencing procedures.** Estonian WGS and Finnish WGS samples were sequenced at the Broad Institute on Illumina HiSeq X Ten machines run to

20 $\times$  and 30 $\times$  mean coverage (150-bp paired reads), respectively. Estonian samples followed a PCR-free sample preparation. Swedish WES and Finnish WES samples were sequenced using either the Agilent SureSelect Human All Exon Kit or the Agilent SureSelect Human All Exon v.2 Kit. Sequencing was performed at the Broad Institute on Illumina GAI, Illumina HiSeq2000 or Illumina HiSeq X Ten. Mean target coverage was 90 $\times$ .

All samples were aligned against the GRCh37 human genome reference and BAM processing was carried out using BWA Picard. Genotype calling was done using GATK Haplotype Caller and was performed at the Broad Institute for all studies.

**Hail software.** To overcome the growing computational challenge of learning from large genomic data sets, we used Hail, an open-source software framework for scalably and flexibly analyzing such data (<https://github.com/broadinstitute/hail>). Hail, under active development, includes support for data import/export, quality control, analysis of population structure and methods for performing both common and rare variant association. Hail is written in Scala (a Java virtual machine language) and builds on open-source software for scalable distributed computing including Hadoop (<http://hadoop.apache.org/>) and Spark (<http://spark.apache.org/>). Hail achieves near-perfect scalability for many tasks and can run on thousands of nodes. Hail automates fault-tolerant distribution of data and computation, greatly simplifying distributed pipeline execution compared to traditional HPC job schedulers like LSF and Grid Engine. Pipelines written in Hail's high-level language typical require orders of magnitude fewer lines of code than comparable pipelines written in general purpose languages.

**Samples and variants quality control.** Quality control was performed independently for each study using Hail. We excluded individuals with high proportion of chimeric reads ( $>5\%$ ), high contamination ( $>5\%$ ) or an excessive number of singletons variants not observed in Exome Aggregation Consortium (ExAC) data ( $> 100$  for WES and  $> 20,000$  for WGS). We included only unrelated individuals (IBD proportion  $< 0.2$ ) and those for whom the sex predicted from genetic data matched the self-reported gender.

We kept only 'PASS' variants, as determined by the Genome Analysis Toolkit<sup>25</sup> Variant Quality Score Recalibration (VQSR) filter set to missing variants with  $GQ < 20$  and allele balance  $> 0.8$  or  $< 0.2$ . We further excluded variants with call rate  $< 0.8$ . In WGS data, we excluded low-complexity regions as defined by Li<sup>26</sup>. In the burden test analysis we excluded variants with both Hardy-Weinberg equilibrium test  $P < 1 \times 10^{-6}$  and negative inbreeding coefficient (expected heterozygosity less than observed heterozygosity).

**Annotation and URV scores definitions.** Annotation was performed using SnpEff 4.2 (build 2015-12-05)<sup>27</sup> using Ensemble gene models from database GRCh37.75. We further annotated variants with SnpSift 4.2 (build 2015-12-05)<sup>28</sup> using annotations from database dbNSFP 2.9 (ref. 29). In **Supplementary Table 3** we provide a detailed description of the criteria used for selecting variants in each score. The set of HC genes was defined separately for disruptive and damaging variants. For disruptive and synonymous mutations we defined HC genes as those having a probability of being loss-of-function intolerant ( $P_{LI}$ )  $> 0.9$  ( $n = 3,488$  genes). For missense damaging mutation we used a missense  $z$ -score  $> 3.09$  ( $n = 1,614$  genes)<sup>30</sup>. Both measures have been previously described<sup>30</sup> and are available online at [ftp://ftp.broadinstitute.org/pub/ExAC\\_release/release0.3/functional\\_gene\\_constraint](ftp://ftp.broadinstitute.org/pub/ExAC_release/release0.3/functional_gene_constraint). We used a version derived from the Exome Aggregation Consortium without cases of psychiatric disorders.

**Principal component analysis and mixed models.** We used a subset of high-confidence single nucleotide polymorphisms (SNPs) to calculate principal components. We selected variants with minor allele frequencies larger than 5%, call rate  $> 90\%$  and Hardy-Weinberg equilibrium test  $P > 1 \times 10^{-6}$ , and we pruned for variants in linkage disequilibrium using Plink with command line '--indep 50 5 2'.

We used a similar approach to filter variants used to generate the genetic relationship matrix (GRM). We then fit a linear mixed model including the GRM as random effect and age, sex, year of birth, (year of birth - 1950)<sup>2</sup>, (year of birth - 1950)<sup>3</sup>, the number of singletons synonymous variants not in ExAC and the number of URVs in HC genes as fixed effects.

**Association between URVs and educational attainment.** We fit a linear regression model where the dependent variable was YOE and the independent predictors were: age, sex, year of birth, (year of birth – 1950)<sup>2</sup>, (year of birth – 1950)<sup>3</sup>, the 10 first principal components, the number of singletons synonymous variants not in ExAC, schizophrenia status (only in studies including schizophrenic patients) and the URV score (count of disruptive, damaging or synonymous URVs). We adjusted for the number of all ultra-rare synonymous variants to correct for potential technical artifacts. We observed similar results when adjusting for the number of ultra-rare synonymous variants + number of ultra-rare disruptive (or damaging) variants in HC genes.

**Brain-expressed and brain-enriched HC genes analysis.** Using Genotype-Tissue Expression Consortia (GTEx) data<sup>14</sup>, we ranked gene-expression levels (in RPKM) in brain tissues and defined the top 20% HC genes as ‘brain-expressed’ ( $n = 683$  genes and 313 genes for disruptive and damaging, respectively). Conversely, we defined ‘non-brain-expressed’ as the bottom 20% of the HC genes ( $n = 683$  genes and 313 genes for disruptive and damaging, respectively).

We also compute estimated fold-change in the brain as follows. Suppose samples 1, 2, ...,  $N_b$  are brain samples and samples ( $N_b + 1$ ), ( $N_b + 2$ ), ...,  $N$  are the samples from other tissues. Denote with  $x_{ij}$  the expression of gene  $j$  and sample  $i$ , in reads per kilobase of transcript per million (RPKM). We compute fold-change (FC):

$$FC_j = \frac{\frac{1}{N_b} \sum_{i=1}^{N_b} x_{ij}}{\frac{1}{N} \sum_{k=1}^N x_{kj}} = \frac{\text{mean}(x_j | \text{brain})}{\text{mean}(x_k)}$$

We labeled the genes  $j$  for which  $FC_j > 2$  as brain-enriched genes and those for which  $FC_j > 0.5$  as non-brain-enriched genes. The numbers of brain-enriched HC genes were 447 and 287 for disruptive and damaging mutations, respectively. The numbers of non-brain-enriched HC genes were 2,225 and 935 for disruptive and damaging mutations, respectively.

**Polygenic score, CNVs and runs of homozygosity.** The polygenic score for YOE was obtained from array data in the Swedish WES study (quality control for the array data have been previously described<sup>20</sup>) and directly from WGS data in the Finnish WGS and Estonian WGS studies. We included all the independent markers with  $P < 1$  in largest the GWAS of educational attainment<sup>5</sup> and obtained the polygenic score as weighted sum of risk alleles using the `--score` command in Plink<sup>31</sup>.

CNVs for the Swedish WES study were called as part of a separate project<sup>32</sup> using a composite pipeline comprising the CNV callers PennCNV, iPattern, Birdsuite and C-Score organized into component pipelines. We considered only rare CNVs by filtering out all CNVs that presented at  $\geq 1\%$  allele frequency. CNVs  $< 20$  kb or having fewer than 10 probes were also excluded. We used the plink `--cnv-intersect` function with a value of 0.5 to determine the overlap between detected CNVs and the list of pathogenic CNVs reported in **Supplementary Table 4**.

CNVs in Finnish WGS and Estonian WGS were genotyped according to the methods described in ref. 33 and implemented in Genome STRiP 2.0. Briefly, read-depth information was collected from WGS data, excluding regions of the genome that were not uniquely alignable or that had low sequence complexity, and adjusted for GC content bias. Each CNV reported in **Supplementary Table 4** was directly genotyped using Genome STRiP’s genotyping module, which examines the read depth across all samples and fits a constrained Gaussian mixture model with components representing each possible diploid copy number and sample-specific variance terms to account for differences in sequencing depth.

The summed runs of homozygosity were determined using the same pipeline described in<sup>15</sup>. Specifically, we used Plink with command line ‘--homozyg --homozyg-window-snp 35 --homozyg-snp 35 --homozyg-kb 1500 --homozyg-gap 1000 --homozyg-density 250 --homozyg-window-missing 5 --homozyg-window-het 1’.

**Gene-based burden test.** We first extracted from each data set variants falling within UCSC known genes and merged the four data sets using Plink. If a variant was not present in all cohorts, we forced it as a homozygous reference across the remaining cohorts (using the ‘--fill-missing-a2’ option in Plink). We then computed principal components for the combined data set after further merging it with 1000 Genomes Project samples, as described in Genovese *et al.*<sup>10</sup>. To test the hypotheses that disruptive URVs in individual genes were associated with YOE and college status, we performed a burden test<sup>34</sup> using the SKAT software<sup>35</sup> using default parameters (method = davies, impute.method = bestguess, r.corr = 1.0), adjusting for age, sex, year of birth, (year of birth – 1950)<sup>2</sup>, (year of birth – 1950)<sup>3</sup>, the first 10 principal components, schizophrenia status and number of URVs identified in coding regions. We used a python wrapper to run the SKAT software (available at <https://github.com/freeseeek/gwaspipeline>).

**Data availability.** Swedish WES data are available through dbGAP at [http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000473.v2.p2](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000473.v2.p2).

The code used in this study is available at: [https://github.com/andgan/URV\\_edu\\_attainment](https://github.com/andgan/URV_edu_attainment).

A **Supplementary Methods Checklist** is available.

20. Ripke, S. *et al. Nat. Genet.* **45**, 1150–1159 (2013).
21. Leitsalu, L. *et al. Int. J. Epidemiol.* **44**, 1137–1147 (2015).
22. Vartiainen, E. *et al. Int. J. Epidemiol.* **39**, 504–518 (2010).
23. Jostins, L. *et al. Nature* **491**, 119–124 (2012).
24. Schizophrenia Working Group of the Psychiatric Genomics Consortium. *Nature* **511**, 421–427 (2014).
25. DePristo, M.A. *et al. Nat. Genet.* **43**, 491–498 (2011).
26. Li, H. *Bioinformatics* **30**, 2843–2851 (2014).
27. Cingolani, P. *et al. Front. Genet.* **3**, 35 (2012).
28. Cingolani, P. *et al. Fly (Austin)* **6**, 80–92 (2012).
29. Liu, X., Jian, X. & Boerwinkle, E. *Hum. Mutat.* **34**, E2393–E2402 (2013).
30. Samocha, K.E. *et al. Nat. Genet.* **46**, 944–950 (2014).
31. Chang, C.C. *et al. Gigascience* **4**, 7 (2015).
32. Marshall, C. *et al. Preprint at bioRxiv* <http://dx.doi.org/10.1101/040493> (2016).
33. Handsaker, R.E. *et al. Nat. Genet.* **47**, 296–303 (2015).
34. Madsen, B.E. & Browning, S.R. *PLoS Genet.* **5**, e1000384 (2009).
35. Ionita-Laza, I., Lee, S., Makarov, V., Buxbaum, J.D. & Lin, X. *Am. J. Hum. Genet.* **92**, 841–853 (2013).