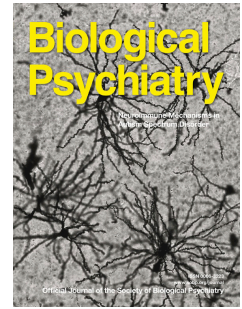


Accepted Manuscript

Developmental changes within the genetic architecture of social communication behaviour: A multivariate study of genetic variance in unrelated individuals

Beate St Pourcain, Lindon J. Eaves, Susan M. Ring, Simon E. Fisher, Sarah Medland, David M. Evans, George Davey Smith



PII: S0006-3223(17)32008-5

DOI: [10.1016/j.biopsych.2017.09.020](https://doi.org/10.1016/j.biopsych.2017.09.020)

Reference: BPS 13333

To appear in: *Biological Psychiatry*

Received Date: 17 September 2016

Revised Date: 1 August 2017

Accepted Date: 17 September 2017

Please cite this article as: St Pourcain B., Eaves L.J, Ring S.M, Fisher S.E, Medland S., Evans D.M & Smith G.D., Developmental changes within the genetic architecture of social communication behaviour: A multivariate study of genetic variance in unrelated individuals, *Biological Psychiatry* (2017), doi: 10.1016/j.biopsych.2017.09.020.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Abstract word length: 248
Article word length: 3998
Tables: 1
Figures: 3
Supplementary Methods: 6
Supplementary Tables: 12
Supplementary Figures: 3

Developmental changes within the genetic architecture of social communication behaviour: A multivariate study of genetic variance in unrelated individuals

Short title: A multivariate study of genetic variance in unrelateds

Beate St Pourcain^{1,2,3#}, Lindon J Eaves⁴, Susan M Ring^{2,5}, Simon E Fisher^{1,4}, Sarah Medland⁶, David M Evans^{2,7}, George Davey Smith^{2,5}

- 1 Language and Genetics Department, Max Planck Institute for Psycholinguistics, The Netherlands
- 2 MRC Integrative Epidemiology Unit (MRC IEU), University of Bristol, UK
- 3 Donders Institute for Brain, Cognition and Behaviour, Radboud University, The Netherlands
- 4 Department of Human and Molecular Genetics, Institute for Psychiatric and Behavioral Genetics, Commonwealth University School of Medicine, Richmond, Virginia, USA
- 5 School of Social and Community Medicine, University of Bristol, UK
- 6 Psychiatric Genetics, QIMR Berghofer Medical Research Institute, Queensland, Australia
- 7 University of Queensland Diamantina Institute, Translational Research Institute, University of Queensland, Queensland, Australia

Corresponding author

Postal address: Max Planck Institute for Psycholinguistics, Wundtlaan 1, 6525 XD Nijmegen, the Netherlands

e-mail: Beate.StPourcain@mpi.nl / Phone: +31 24 3521964 / Fax: +31 24 3521213

Abstract

Background: Recent analyses of trait-disorder overlap suggest that psychiatric dimensions may relate to distinct sets of genes that exert their maximum influence during different periods of development. This includes analyses of social-communication difficulties that share, depending on their developmental stage, stronger genetic links with either Autism Spectrum Disorder or schizophrenia. Here we developed a multivariate analysis framework in unrelated individuals to model directly the developmental profile of genetic influences contributing to complex traits, such as social-communication difficulties, during a ~10-year period spanning childhood and adolescence.

Methods: Longitudinally assessed quantitative social-communication problems ($N \leq 5,551$) were studied in participants from a UK birth cohort (ALSPAC, 8 to 17 years). Using standardised measures, genetic architectures were investigated with novel multivariate genetic-relationship-matrix structural equation models (GSEM) incorporating whole-genome genotyping information. Analogous to twin research, GSEM included Cholesky decomposition, common pathway and independent pathway models.

Results: A 2-factor Cholesky decomposition model described the data best. One genetic factor was common to SCDC measures across development, the other accounted for independent variation at 11 years and later, consistent with distinct developmental profiles in trait-disorder overlap. Importantly, genetic factors operating at 8 years explained only ~50% of the genetic variation at 17 years.

Conclusion: Using latent factor models, we identified developmental changes in the genetic architecture of social-communication difficulties that enhance the understanding of ASD and schizophrenia-related dimensions. More generally, GSEM present a framework for modelling shared genetic aetiologies between phenotypes and can provide prior information with respect to patterns and continuity of trait-disorder overlap.

Keywords: ALSPAC, Structural equation modelling, Longitudinal analysis, Genetic variance decomposition, Genetic-relationship matrix structural equation modelling, Genetic relationship matrix

Introduction

The extent to which genetic aetiologies are shared between traits and disorders naturally depends on the genetic composition of the two phenotypes. While psychiatric disorders are diagnostic entities, defined by clinical criteria including the age of onset, human behaviour changes continuously during development. This includes developmental alterations in complex genetic trait architectures as reported for cognitive (1) but also social-communication related characteristics (2).

Difficulties to socially engage and communicate with others, as observed in the general population, are heritable (twin- $h^2=0.74$) (3) and a considerable proportion of the underlying genetic variation can be tagged by Single Nucleotide Polymorphisms (SNPs, $SNP-h^2 \leq 0.45$) (2). For both, social-communication and social interaction problems, multivariate twin (4;5) and bivariate GREML (genetic-relationship-matrix residual maximum likelihood) studies (6) reported evidence for a degree of genetic stability, but also change during childhood and adolescence (2;7;8) that may affect genetic similarities with other traits.

Studying the genetic overlap between psychiatric illness and social-communication difficulties across multiple developmental stages, different developmental profiles for childhood- versus adult-onset psychiatric disorders have been identified (9). The genetic overlap with clinical Autism Spectrum Disorder (ASD), a complex highly heritable early-onset neurodevelopmental condition (10), was strongest for social-communication difficulties during childhood, but declined with progressing age of the trait. By contrast, the genetic correlation with clinical schizophrenia, an adult-onset psychiatric illness with a typical first-time diagnosis between 16 to 30 years (10), was highest for social-communication problems during later adolescence (9). Thus, the risk of developing these contrasting psychiatric conditions might be related to distinct sets of genes, both of which affect social

communication skills, but exert their maximum influence during different periods of development.

Discontinuity in trait-disorder overlap may, however, also result because of attrition-related artefacts such as decreasing power or inherent sample bias (11). As knowledge about developmental changes in complex genetic trait architectures is still scarce, development-related variations in trait-disorder overlap are often dismissed.

The aim of this work is to provide insight into the developmental profile of genetic factors influencing complex traits, such as social-communication difficulties during childhood and adolescence, using a longitudinal analysis framework. Building on our previous work (2;9), we investigate here two extreme hypotheses: We evaluate whether the genetic variance/covariance structure of social-communication difficulties during childhood and adolescence is consistent with multiple independent genetic influences, suggesting developmental changes in the genes responsible for inter-individual variation over time, or whether, alternatively, there is evidence for a shared single genetic factor, irrespective of age.

To study the developmental profile of genetic factors in unrelated individuals, we implemented multivariate genetic-relationship-matrix structural equation models (GSEM). These models utilise genome-wide genetic relationship matrices (GRMs)(12), calculated from hundreds of thousands of SNPs across the genome, to estimate the total amount of phenotypic variance and covariance tagged by common genetic variants, similar to GREML (12;13). GREML and related approaches (12;14–16) have re-shaped the research of complex genetic trait architectures beyond twin designs by exploiting the availability of genome-wide genetic data in cohorts of unrelated individuals. Genetic correlations are, however, typically estimated by these methods by studying two phenotypes only. Using a structural equation modelling (SEM) framework (17), as widely applied within twin research (4;5), we now

extend this bivariate approach by flexibly modelling complex latent genetic factor structures within a multivariate context.

In this paper we use multivariate GSEM to model longitudinal data on social-communication difficulties across childhood and adolescence in the Avon Longitudinal Study of Parents and Children (ALSPAC), a phenotypically-rich longitudinal population-based birth cohort from the UK (18).

Methods

Participants and measures

All analyses were carried out using children's data from ALSPAC, a UK population-based longitudinal pregnancy-ascertained birth-cohort (estimated birth date: 1991 to 1992)(18). Please note that the study website contains details of all the data that is available through a fully searchable data dictionary (<http://www.bris.ac.uk/alspac/researchers/data-access/data-dictionary/>). Ethical approval was obtained from the ALSPAC Law-and-Ethics Committee (IRB00003312) and the Local Research-Ethics Committees. Written informed consent was obtained from a parent or individual with parental responsibility and assent (and for older children consent) was obtained from the child participants.

Phenotype information: Social-communication difficulties during childhood and adolescence were collected with the 12-item mother-reported Social Communication Disorder Checklist (SCDC; score-range: 0 to 24, age range: 3 to 18 years)(3). The SCDC is a brief screening instrument of social reciprocity and verbal/nonverbal communication (e.g. "Not aware of other people's feelings"), which has high reliability and internal consistency, and good validity (3) with higher scores reflecting more social-communication deficits. Quantitative SCDC scores in ALSPAC children and adolescents were repeatedly measured at

8, 11, 14 and 17 years and information on phenotypic and genotypic data was available for 4,174 to 5,551 children (Supplementary Table S1).

Descriptive analyses of SCDC scores were carried out in R.v.3.2.4. The distribution of SCDC scores was positively skewed and predominantly leptokurtic (Supplementary Table S1). Each score was adjusted for sex, age and the two most significant ancestry-informative principal components (see below) using ordinary least square (OLS) regression. Residuals were subsequently transformed to perfect normality using rank-based inverse normal transformation (19), as previously reported (9), to allow for comparisons across different algorithms (see below). There were moderate phenotypic correlations between repeatedly assessed SCDC scores, using both untransformed and transformed data (Supplementary Table S2, SCDC: Spearman's- ρ : 0.39 to 0.57; Pearson- r : 0.38 to 0.61) as previously shown (9).

Genome-wide genotype information: ALSPAC children were genotyped using the Illumina HumanHap550 quad chip genotyping platforms (Supplementary Methods). After quality control, 8,237 children and 477,482 directly genotyped Single Nucleotide Polymorphisms (SNPs) were kept within the study.

GSEM

Multivariate SEM techniques were used to assess the relative importance of genetic and residual influences to variation in longitudinal SCDC scores during child and adolescent development. Similar to GREML (12), GSEM use the genetic similarity between unrelated individuals to partition the expected phenotypic variance/covariance matrix into genetic and residual components. More generally, however, the statistical framework of GSEM is analogous to twin analysis methodologies (4;5), but uses GRMs, instead of twin correlations, to estimate genetic variance/covariance structures using full information maximum likelihood

(FIML). Thus, genetic and environmental influences are modelled in the GSEM framework as latent factors contributing to inter-individual covariation in phenotypic measures. The advantage of our approach is that multivariate SEM methodology has been widely established within twin research (4;5) and allows for flexible modelling of complex genetic factor structures. Conversely, GREML, as implemented in the GCTA software package, is currently restricted to bivariate situations (20). While multivariate GSEM can be fit with SEM software such as OpenMx (21) using both mxGREML and FIML algorithms, these models are currently computationally expensive (see Results). We therefore implemented GSEM within R (Rv3.2.4) (for details see Supplementary Methods).

In short, GSEM describe the phenotypic covariance structure using one or more additive genetic factors A that capture genetic variance, tagged by common genotyped SNPs, as well as one or more residual factors E that capture residual variance, containing untagged genetic variation and unique environmental influences (including measurement error). As SEM methodology has its origins in the method of path analysis (22), path diagrams are useful in visualising the relationship among observed and latent variables (represented as squares and circles respectively, see e.g. Figure 2). Single headed arrows (factor loadings or 'paths') denote causal relationships between measures, whereas double headed arrows define correlations.

In our formulation, additive genetic variances (GSEM-Var_g) and genetic covariances (GSEM-Cov_g) are modelled as the product of additive genetic factor loadings and genetic factor variances (the latter being standardised to unit variance). For example, using multivariate GSEM, a saturated model can be fit to the data through a decomposition of both the genetic variance and residual variance into as many latent factors as there are observed variables (Cholesky decomposition model; see Supplemental methods). Estimated genetic variances and covariances can then be used to estimate genetic correlations (GSEM-r_g) (23),

i.e. the extent to which two phenotypes share common genetic factors (Supplementary Methods). Here, we utilised the Cholesky decomposition model as saturated and baseline model (Supplementary Information). Beside Cholesky decomposition models, multivariate GSEM also permit the fitting of models with smaller numbers of latent genetic and residual factors, defined according to theory (24).

Multivariate GSEM of longitudinally assessed SCDC scores were fitted in two stages.

In a first step (I), we specified *a priori* three standard multivariate AE models, analogous to twin research (Figure 2A-C): we studied a Cholesky decomposition model (saturated model), an independent pathway model and a common pathway model.

- 1) The Cholesky decomposition model, as described above, is a fully parametrised descriptive model without any restrictions on the structure of latent genetic and residual influences (20 free parameters) (Figure 2A) and involves multiple independent genetic influences sharing genetic aetiologies across development.
- 2) The independent pathway model, in its simplest form, specifies a single common genetic factor and a single common residual factor, in addition to age-specific genetic and residual influences (16 free parameters) (Figure 2B).
- 3) The common pathway model, in its simplest form, parametrises a single latent factor, influenced by both genetic and residual sources of variance, in addition to age-specific genetic and residual influences (Figure 2C), and is the most constrained model (14 free parameters). The model constrains the variance of the latent factor to one (i.e. the sum of squared genetic and residual factor loadings). Although the likelihood of this model can be estimated, the resulting Hessian is not invertible due to singularity problems. For these reasons, the model constraint was relaxed within this work.

Both, the independent pathway model and the common pathway model are consistent with a shared single genetic factor across development and are nested submodels of the full Cholesky decomposition model.

The goodness-of-fit of GSEM to empirical data was assessed using likelihood ratio test (LRT), the Akaike Information Criterion (AIC) (25) and the Bayesian Information Criterion (BIC) (26) (Supplementary Methods).

In a second step (II), we adopted a data-driven approach and investigated the pattern of genetic factor loadings for the best fitting model from (I) in detail. The smallest genetic factor loadings were successively dropped from the model and the overall fit of the model compared with the best-fitting *a priori* defined GSEM (or an adapted form) using LRTs. The statistical significance of factor loadings was assessed using a Wald test (2-sided test). Standard errors (SEs) for genetic and residual variances and covariances, and genetic correlations were derived from the variance-covariance matrix of the estimated factor loadings using the delta method. Standard errors for factor loadings were estimated by GSEM. Note that for rank-transformed measures with unit variance, such as the SCDC scores in this study, genetic variances are equivalent to $\text{SNP-}h^2$ estimates. However, path coefficients for multivariate GSEM were re-standardised to enhance the interpretability.

GRMs were estimated using the GCTA software (12) and based on directly genotyped SNPs. All GSEM were fitted to data from participants with non-missing information to simplify the estimation algorithm. All R scripts are available via the R gsem package. (<https://gitlab.gwdg.de/beate.stpourcain/gsem>, Supplementary Information).

For the purpose of benchmark comparisons with univariate GCTA, we also fitted univariate GSEM, where genetic variances were estimated as a single variance component.

GREML

The GCTA software package can be used to estimate the proportion of phenotypic variation that is jointly explained by SNPs on a genotyping chip using GREML (13) (AE model). Likewise, bivariate GREML (20) allows estimating genetic covariances and genetic correlations between two phenotypes. An advantage of this method is that genetic correlations between two phenotypes can be estimated even when these phenotypes are not measured in the same individuals.

Univariate and bivariate GREML were carried out as part of sensitivity and simulation analyses. For comparison with GSEM, genetic relationship matrices (GRMs) were derived from directly genotyped SNPs, but excluded individuals with a pairwise relationship >0.025 , as recommended (13). All analyses were conducted with GCTA software v1.25.2 (12).

OpenMx SEM models

OpenMx SEM models (21), as implemented in the OpenMx software (<http://openmx.psyc.virginia.edu/>)(v2.5 and v2.7), were fitted using FIML and mxGREML and included a full Cholesky decomposition of both genetic and residual variances (AE model, see above). Bivariate OpenMx SEM analyses were conducted as part of a simulation analysis. Genetic variances, genetic covariances, and genetic correlations were derived as described for GSEM above.

All analyses were conducted on High Performance Clusters at the University of Bristol and the MPI for Psycholinguistics.

Data simulation

To evaluate the accuracy of multivariate GSEM, we carried out data simulations (Supplementary Methods).

Attrition analysis

SCDC-attrition scores were generated to investigate potential sources of bias. Analyses included sample-specific estimates of genetic correlations among SCDC-attrition scores, and between SCDC scores and subsequent sample dropout (Supplementary Methods).

Results

Accuracy of multivariate GSEM

We simulated a bivariate trait (N=5000) with two standardised measures (10 replicates; Supplementary Figure S1A, Supplementary Table S3) and confirmed the accuracy of multivariate GSEM through comparison with GCTA and OpenMx software. All methods provided accurate estimates, both with respect to genetic and residual variances and covariances as well as genetic and residual factor loadings (GSEM and OpenMx SEM models only), with comparable RMSE, MAD and little bias ($\text{Bias}^2 < 10^{-3}$ for all methods, Supplementary Table S3). Computationally, multivariate OpenMx SEM models were, however, more expensive (≤ 78 GB RAM FIML v2.5; ≤ 2694 minutes mxGREML/FIML v2.7) than multivariate GSEM (≤ 13 GB RAM, ≤ 301 minutes) per single bivariate replicate analysis. A comparison of computing resources is shown in Supplementary Table S4. There was also little difference between estimated OpenMx versus GSEM parameters when analysing a trivariate simulated trait with three standardised measures, as part of a benchmark

test (Supplementary Figure S1B, Supplementary Table S5). Note that trivariate replicate analyses using OpenMx were not considered within this study due to computational constraints.

Univariate analyses

Using univariate GSEM, common genetic variants explained a large proportion of phenotypic variation in SCDC scores during childhood as well as during later adolescence (age 8: $\text{Var}_g(\text{SE})=0.25(0.061)$, $p=3.4\times 10^{-5}$; age 11: $\text{Var}_g(\text{SE})=0.22(0.061)$, $p=2.9\times 10^{-4}$, age 17: $\text{Var}_g(\text{SE})=0.47(0.086)$, $p=4.4\times 10^{-8}$; Figure 1, Supplementary Table S6) but not during early adolescence (age 14, $\text{Var}_g(\text{SE})=0.086(0.064)$, $p=0.18$), as previously reported (2). Univariate GCTA(GREML) yielded nearly identical results (Supplementary Table S7).

Figure 1 about here

Multivariate analyses

We first examined the profile of genetic factors contributing to variation in SCDC scores during development (13,180 observations; 3,295 participants) using three *a priori* defined multivariate GSEM (Figure 2A-C). Based on all three fit indices, LRT, AIC and BIC, the best-fitting *a priori* defined model was the full Cholesky decomposition model (Model 1, Table 1, Figure 2A, Figure 3A). Neither a single factor independent pathway model nor a single factor common pathway model could sufficiently capture the underlying variance/covariance structure of the data. As the full Cholesky decomposition model is, however, also the baseline model, the model identification progressed with the identification

of meaningful GSEM through data-driven model modifications. Consistent with near zero factor loadings for the latent genetic factors A_3 and A_4 (Supplementary Table S8), a two genetic factor Cholesky model was studied (Model 4, Figure 2D) that provided a near-identical fit to the data (Table 1, $\Delta X^2 < 0.01$ ($\Delta df=3$), $p=1$). This model parametrised one genetic factor arising at age 8 years, and a second independent genetic factor explaining novel genetic influences arising at age 11 years, each contributing to phenotypic variation during later development (Figure 2D). Using LRTs, the model fitting progressed (Model 5, Table 1, Supplementary Table S8) until all genetic factor loadings reached $p < 0.05$ without a significant drop in the log-likelihood ($\Delta X^2 = < 0.01$ ($\Delta df=2$), $p=1$, with respect to Model 4).

Figure 2 about here

Table 1 about here

The identified model included one common genetic factor A_1 , accounting for shared phenotypic variation throughout development, as well as a second genetic factor A_2 influencing SCDC scores at 11 years and especially at 17 years of age (Table 1, Figure 3B). Figure 3 shows the full Cholesky decomposition model (Model 1) and its best-fitting reduced form (Model 5) with their standardised path coefficients (factor loadings ≥ 0.32 explain $>10\%$ of the phenotypic variance).

Overall, the estimates of genetic variance, as predicted by GSEM (Model 1 and 5, Supplementary Table S9), were consistent with univariate GSEM estimates (Figure 1), although latter were based on larger sample numbers (Supplementary Table S6). The pattern of genetic factor loadings suggested, however, a dynamic change in the variance composition

of the trait during development such that only ~50% of the genetic variance at age 17 was accounted for by genetic variation at age 8 (e.g. age 17: ratio $\text{Var}_g(A_1)$ to $\text{Var}_g(A_1+A_2)$; Model 1: 0.53(SE=0.18)%; Model 5: 0.53(SE=0.12%)(Figure 1).

Figure 3 about here

The predicted bivariate genetic correlations by multivariate GSEM (Model 1 and 5, Supplementary Table S9) were overall similar to bivariate GCTA(GREML) estimates, although latter were based on larger numbers of observations (Supplementary Table S10 and Supplementary Figure S3). Restricting analyses to the same sets of individuals, both bivariate GSEM and bivariate GCTA(GREML) provided near-identical estimates (Supplementary Table S10), although these analyses were less powerful. Thus, small differences in genetic correlations patterns, as estimated by multivariate GSEM versus bivariate GCTA(GREML), are likely to be due to minor differences in sample numbers.

There was furthermore little evidence that genetic influences between SCDC scores and subsequent SCDC sample dropout are shared in ALSPAC (Supplementary Table S11). Nominal evidence for a genetic correlation was observed between SCDC scores at 8 years and dropout at 14 years only ($r_g=0.39(\text{SE}=0.19)$, $p_{\text{one-tailed}}=0.02$). Nonetheless, SCDC attrition scores were genetically correlated across all SCDC measures in ALSPAC ($p_{\text{one-tailed}}<10^{-3}$, Supplementary Table S12).

Discussion

Using multivariate SEM in combination with common variant-based genetic correlation matrices, we investigated the developmental structure of genetic factors contributing to social-communication difficulties during childhood and adolescence. We showed that the genetic architecture of this population-based complex trait changes continuously during development and is consistent with multiple genetic influences operating at different stages during development. Thus, our study provides evidence against the hypothesis that social communication behaviour during development is a genetically homogenous phenotype.

The best-fitting model, specifying two distinct genetic factors, suggested that the genetic origins of child and adolescent social-communication behaviour lie in middle and late childhood. The first genetic factor, parametrised to account for all genetic influences at age 8 years, explained a considerable proportion of phenotypic variance throughout development (>20%) with the exclusion of SCDC scores at age 14 that have negligible SNP- h^2 estimates. (This is consistent with recent reports of low SNP- h^2 for autistic symptoms at the beginning of adolescence (1) and might be related to pubertal adjustments (2)).

The second genetic factor, parametrised to be independent of the first one and to capture novel genetic influences arising at age 11 years, explained predominantly phenotypic variation at 17 years of age (~19%). Thus, the model predicted changes in the composition of the genetic variance during development, and only ~50% of the genetic variation at age 17 was accounted for by genetic variation at age 8. Within defined developmental stages, however, such as those spanning mid-childhood to very early adolescence (e.g. 8 to 11 years), we found evidence for strong genetic correlations across measures. These results are consistent with recent longitudinal twin research that reported moderate to high genetic stability for autistic traits, including communication impairments, between mid-childhood

and early adolescence (7), but only moderate genetic stability between behaviour in childhood versus emerging adulthood (8). The identified genetic factor structure using GSEM reflects therefore both a degree of genetic stability, but also genetic change in social-communication behaviour during development, depending on the size of the developmental window.

The identification of two distinct genetic factors, especially during later adolescence, suggests that SCDC scores at age 8 or 11 years are, in terms of average composition, different from those influencing SCDC scores at age 17. Developmental changes in the genetic architecture of social communication traits are consistent with biological maturation processes during childhood and adolescence. For example, synaptic pruning in the cerebral cortex is a signature late maturational process for generating a diversity of neuronal connections (27), which occurs during puberty and extends into early adult life (28). In parallel, there are changes in adolescent social cognitive development, especially with respect to emotional perspective taking, resistance to peer influence and changes in social behaviour (29). Given the identified genetic factor structure, it could be speculated whether multiple concepts of 'social reciprocity and verbal/nonverbal communication' may co-exist, especially at age 17, and whether changes in genetic factor contributions may continue into early adulthood. Thus, even for psychological instruments with high reliability, internal consistency and good discriminant validity, like the SCDC (3), the nature of the captured continuous phenotype may vary across developmental periods spanning ~10 years. This underlines the need for behavioural genetic studies across the life-span.

An important implication that flows from the observation of developmental variations in the genetic trait architecture is that measures assessed at different developmental stages may reveal different patterns of trait-disorder overlap, as previously shown for clinical ASD and schizophrenia respectively (9). Moreover, the identification of a 2-genetic factor is also

consistent with recent reports of little genetic overlap between ASD versus schizophrenia-related dimensions (30), especially with respect to social-communication symptoms. Structural models capturing developmental changes in the genetic architecture of complex phenotypes can therefore be leveraged to obtain prior information concerning the stability of trait-disorder overlap and consequently the extent to which development-specific genetic trait factors are shared among different psychiatric dimensions.

Our findings have therefore specific relevance for the study of functional dimensions of human behaviour spanning the continua from normal to abnormal and across development, consistent with the framework of Research Domain Criteria (31).

Finally, our study proves that structural models of genetic influences in unrelated individuals, as captured by GRMs, are computationally feasible within a longitudinal context. Beyond the scope of bivariate GCTA(GREML), multivariate GSEM allow for the modelling of complex latent genetic factor structures across different stages of development, in particular their genetic variance composition, and can reveal developmental origins of genetic variation that are otherwise hidden. It is furthermore possible to envisage that the concept of GSEM can be extended to investigate multivariate models of cross-disorder overlap and other complex phenomena, such as reciprocal causation. Note that also novel OpenMx FIML and mxGREML algorithms are currently being developed.

A limitation of our study is the analysis of non-missing data across all repeatedly assessed measures. Thus, weaker genetic links, spanning wider age gaps, may not have been sufficiently captured as a consequence of lower power, although genetic correlations predicted by multivariate GSEM and bivariate GCTA(GREML) were overall similar. In addition, cohort studies can be affected by attrition bias (32). We identified, however, little evidence for a specific genetic link between variation in SCDC scores and subsequent sample dropout, although attrition scores across all assessed SCDC measures were genetically

correlated. This is consistent with studies reporting an association between study non-participation, including SCDC dropout, and polygenic risk for schizophrenia (9;32), irrespective of when phenotypes were sampled during development. In addition, we exclusively studied rank-transformed phenotypes to ensure multivariate normality and comparability across different estimation algorithms, and we can therefore not exclude transformation-related biases. However, genetic overlap with psychiatric conditions provided some evidence for the content validity of the analysed trait (9). Also, maternal characteristics may have contributed to phenotypic and, to a lesser extent, genetic correlations. However, the impact of these effects is likely to be small, given the identified developmental changes in genetic variances and covariances for SCDC scores during development. Finally, a Cholesky decomposition of a variance/covariance matrix may not always result in fitting statistics that follow the expected chi-squared distribution (33). Model comparisons using real and simulated data provided, however, little evidence for systematic differences between GCTA(GREML), GSEM and OpenMx SEMs. Thus, despite potential limitations, our study demonstrates that structural models of longitudinally assessed behavioural traits can inform on developmental changes in genetic trait architectures as tagged by common SNPs.

Conclusions

The genetic architecture of social-communication difficulties, as tagged by common genetic variation, changes with age and involves multiple genetic factors operating at different developmental stages during a 10-year period spanning childhood and adolescence. The identification of distinct genetic trait factors is consistent with different profiles of trait-disorder overlap, and underlines the importance of investigating genetic trait variances within a multivariate context.

Acknowledgments

We are extremely grateful to all the families who took part in this study, the midwives for their help in recruiting them, and the whole ALSPAC team, which includes interviewers, computer and laboratory technicians, clerical workers, research scientists, volunteers, managers, receptionists and nurses. This publication is the work of the authors and they will serve as guarantors for the contents of this paper. The UK Medical Research Council and the Wellcome Trust (102215/2/13/2) and the University of Bristol provide core support for ALSPAC. The ALSPAC GWAS data was generated by Sample Logistics and Genotyping Facilities at the Wellcome Trust Sanger Institute and LabCorp (Laboratory Corporation of America) using financial support from 23andMe. Autism Speaks (7132) provided support for autistic-trait related analyses in ALSPAC to BSP. BSP and SF are supported by the Max Planck Society. We thank Robert Kirkpatrick for helpful discussions on structural equation models and support with the OpenMx code. We thank Gregory Carey for his contribution and many helpful discussions as part of the initial work on bivariate mGCTA models carried out together with LE, DE and BSP. We thank Callum Wright and Tobias van Valkenhoef for their help with the High Performance Computing systems.

Conflict of interest

The authors report no biomedical financial interests or potential conflict of interests.

References

1. Trzaskowski M, Yang J, Visscher PM, Plomin R (2014): DNA evidence for strong genetic stability and increasing heritability of intelligence from age 7 to 12. *Mol. Psychiatry* 19: 380–384.
2. St Pourcain B, Skuse DH, Mandy WP, Wang K, Hakonarson H, Timpson NJ, et al. (2014): Variability in the common genetic architecture of social-communication spectrum phenotypes during childhood and adolescence. *Mol. Autism* 5: 18.

3. Skuse DH, Mandy WPL, Scourfield J (2005): Measuring autistic traits: heritability, reliability and validity of the Social and Communication Disorders Checklist. *Br. J. Psychiatry* 187: 568–572.
4. Neale M, Maes HHM (2004): *Methodology for genetic studies of twins and families.*, Dordrecht: Kluwer Academic Publishers.
5. Martin NG, Eaves LJ (1977): The genetical analysis of covariance structure. *Heredity* 38: 79–95.
6. Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR (2012): Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinformatics* 28: 2540–2542.
7. Holmboe K, Rijdsdijk FV, Hallett V, Happé F, Plomin R, Ronald A (2014): Strong Genetic Influences on the Stability of Autistic Traits in Childhood. *J. Am. Acad. Child Adolesc. Psychiatry* 53: 221–230.
8. Taylor MJ, Gillberg C, Lichtenstein P, Lundström S (2017): Etiological influences on the stability of autistic traits from childhood to early adulthood: evidence from a twin study. *Mol. Autism* 8: 5.
9. St Pourcain B, Robinson EB, Anttila V, Bulik-Sullivan B, Maller JB, Golding J, et al. (2017): ASD and schizophrenia show distinct developmental profiles in common genetic overlap with population-based social-communication difficulties. *Mol. Psychiatry* Advance online publication: 10.1038/mp.2016.198.
10. American Psychiatric Association (1994): *Diagnostic and Statistical Manual of Mental Disorders, 4th ed.* Washington, DC: American Psychiatric Association.
11. Martin J, Tilling K, Hubbard L, Stergiakouli E, Thapar A, Smith GD, et al. (2016): Association of Genetic Risk for Schizophrenia With Nonparticipation Over Time in a Population-Based Cohort Study. *Am. J. Epidemiol.* 10.1093/aje/kww009.
12. Yang J, Lee SH, Goddard ME, Visscher PM (2011): GCTA: A Tool for Genome-wide Complex Trait Analysis. *Am. J. Hum. Genet.* 88: 76–82.
13. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. (2010): Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* 42: 565–569.
14. Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Schizophrenia Working Group of the Psychiatric Genomics Consortium, et al. (2015): LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* 47: 291–295.
15. Loh P-R, Tucker G, Bulik-Sullivan BK, Vilhjálmsson BJ, Finucane HK, Salem RM, et al. (2015): Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat. Genet.* 47: 284–290.
16. Speed D, Hemani G, Johnson MR, Balding DJ (2012): Improved Heritability Estimation from Genome-wide SNPs. *Am. J. Hum. Genet.* 91: 1011–1021.
17. Bollen KA (1989): *Structural Equations with Latent Variables*, 1 edition. New York: Wiley-Blackwell.
18. Boyd A, Golding J, Macleod J, Lawlor DA, Fraser A, Henderson J, et al. (2013): Cohort Profile: The “Children of the 90s”—the Index Offspring of the Avon Longitudinal Study of Parents and Children. *Int. J. Epidemiol.* 42: 111–27.
19. Peng B, Yu RK, DeHoff KL, Amos CI (2007): Normalizing a large number of quantitative traits using empirical normal quantile transformation. *BMC Proc.* 1: S156.
20. Lee SH, Yang J, Goddard ME, Visscher PM, Wray NR (2012): Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived

- genomic relationships and restricted maximum likelihood. *Bioinformatics* 28: 2540–2542.
21. Boker S, Neale M, Maes H, Wilde M, Spiegel M, Brick T, et al. (2011): OpenMx: An Open Source Extended Structural Equation Modeling Framework. *Psychometrika* 76: 306–317.
 22. Wright S (1921): Correlation and causation. *J. Agric. Res.* 20: 557–585.
 23. Falconer PDS, Mackay PTF (1995): *Introduction to Quantitative Genetics*, 4 edition. Essex, England: Longman.
 24. MacCallum RC, Roznowski M, Necowitz LB (1992): Model modifications in covariance structure analysis: The problem of capitalization on chance. *Psychol. Bull.* 111: 490–504.
 25. Akaike H (1987): Factor analysis and AIC. *Psychometrika* 52: 317–332.
 26. Schwarz G (1978): Estimating the Dimension of a Model. *Ann. Stat.* 6: 461–464.
 27. Selemon LD (2013): A role for synaptic plasticity in the adolescent development of executive function. *Transl. Psychiatry* 3: e238.
 28. Petanjek Z, Judaš M, Šimić G, Rašin MR, Uylings HBM, Rakic P, et al. (2011): Extraordinary neoteny of synaptic spines in the human prefrontal cortex. *Proc. Natl. Acad. Sci. U. S. A.* 108: 13281–13286.
 29. Burnett S, Blakemore S-J (2009): The Development of Adolescent Social Cognition. *Ann. N. Y. Acad. Sci.* 1167: 51–56.
 30. Taylor MJ, Robinson EB, Happé F, Bolton P, Freeman D, Ronald A (2015): A longitudinal twin study of the association between childhood autistic traits and psychotic experiences in adolescence. *Mol. Autism* 6: .
 31. Cuthbert BN, Insel TR (2013): Toward the future of psychiatric diagnosis: the seven pillars of RDoC. *BMC Med.* 11: 126.
 32. Martin J, Tilling K, Hubbard L, Stergiakouli E, Thapar A, Smith GD, et al. (2016): Association of Genetic Risk for Schizophrenia With Nonparticipation Over Time in a Population-Based Cohort Study. *Am. J. Epidemiol.* kww009.
 33. Carey G (2005): Cholesky Problems. *Behav. Genet.* 35: 653–665.

Table 1: Multivariate GSEM of SCDC scores

Model	Path diagram	-2LL	k	ΔX^2 to model 1	Δdf to model 1	<i>p</i>	AIC	BIC
<i>A priori</i> defined multivariate GSEM								
1. Full Cholesky decomposition model - <i>saturated model</i>	Figure 2A, Figure 3A	7900.97	20	-	-	-	7940.97	8062.97
2. Independent pathway model	Figure 2B	7914.51	16	13.55	4	0.0089	7946.51	8044.12
3. Common pathway model	Figure 2C	8082.7	14	181.73	6	<10 ⁻¹⁵	8110.70	8196.10
Data-driven model modification								
4. Two genetic factor Cholesky model	Figure 2D	7900.96	17	<0.01	3	1	7934.96	8038.67
Best-fitting model								
5. Two genetic factor Cholesky model (excluding non-significant paths) *	Figure 3B	7900.96	15	<0.01	2	1	7930.96	8022.47

The goodness-of-fit of genetic-relationship-matrix structural equation models (GSEM) was assessed with likelihood ratio tests, the Akaike Information Criterion (AIC) and the Bayesian Information criterion (BIC). Following the investigation of *a priori* defined GSEM, the model fitting progressed until all genetic factor loadings reached $p < 0.05$ without a significant drop in the log-likelihood. Path diagrams are shown in Figure 2. The best-fitting model (*) is starred. 3,295 participants had SCDC scores across all ages. k - Number of parameters; LL - Log-likelihood; SCDC - Social and Communication Disorders Checklist

Figures legends**Figure 1:** Genetic variance of SCDC scores during development

Genetic variances for SCDC scores across development as estimated using a univariate model (Supplementary Table S6, $N \geq 4,174$) and the full Cholesky decomposition model (Table 1, Model 1; Supplementary Table S8, $N=3,295$). Genetic factor A_3 and A_4 of the Cholesky decomposition model are not shown as their estimated Var_g was negligible (<0.01). All reported Var_g estimates are equivalent to $\text{SNP-}h^2$ estimates. Grey lines indicate one standard error (SE) in total genetic variance (Var_g) for each SCDC measure.

SCDC - Social and Communication Disorders Checklist; Var_g - Genetic variance

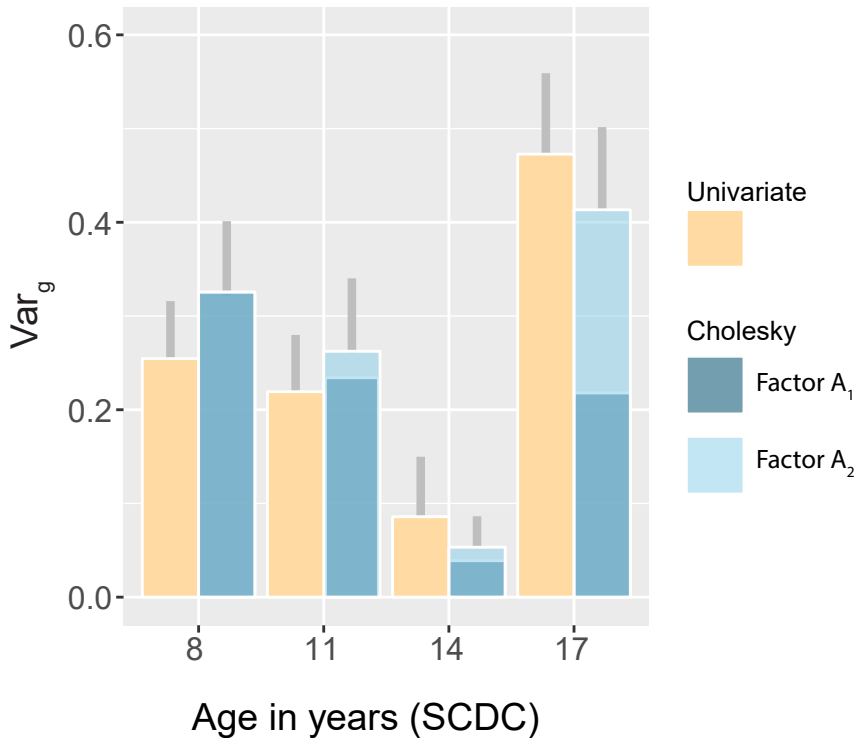
Figure 2: Path diagrams of *a priori* defined multivariate GSEM and data-driven model modifications

A - Full Cholesky decomposition model; B - Independent pathway model; C - Common pathway model; D - Two genetic factor Cholesky model (Data-driven model modification)

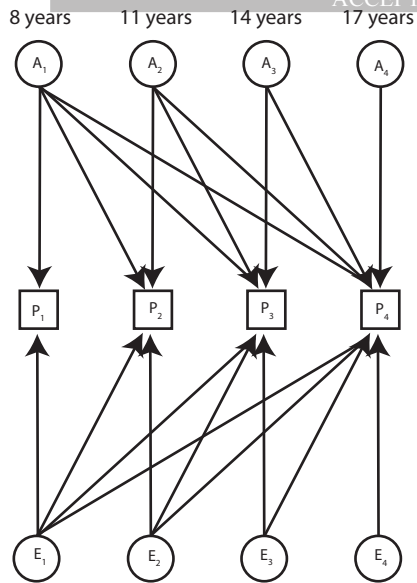
Observed phenotypic measures are represented by squares and latent factors by circles. Single headed arrows ('paths') define causal relationships between variables. Double headed arrows define correlations. Note that the variance of latent variables is constrained to unit variance, this is omitted from the diagrams to improve clarity. GSEM - Genetic-relationship-matrix structural equation models

Figure 3: Path diagram of the full Cholesky decomposition model for SCDC scores (A) and its reduced form (B)

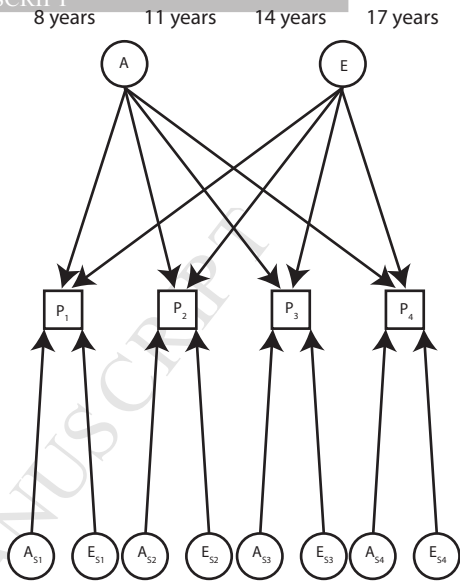
The full Cholesky decomposition model (A) and its most parsimonious reduced form (B) are described in detail in Table 1 (Model 1 and 5 respectively). Corresponding to the phenotypic measures P_1 (8 years), P_2 (11 years), P_3 (14 years) and P_4 (17 years), the latent genetic factors with factor loadings a are A_1 (8 years), factor A_2 (11 years), factor A_3 (14 years), factor A_4 (17 years) and the latent residual factors with factor loadings e are E_1 (8 years), factor E_2 (11 years), factor E_3 (14 years), factor E_4 (17 years). All path coefficients are standardised. 3,295 participants had repeated scores across all ages. Note that the variance of latent variables is constrained to unit variance, this is omitted from the diagrams to improve clarity. SCDC - Social and Communication Disorders Checklist



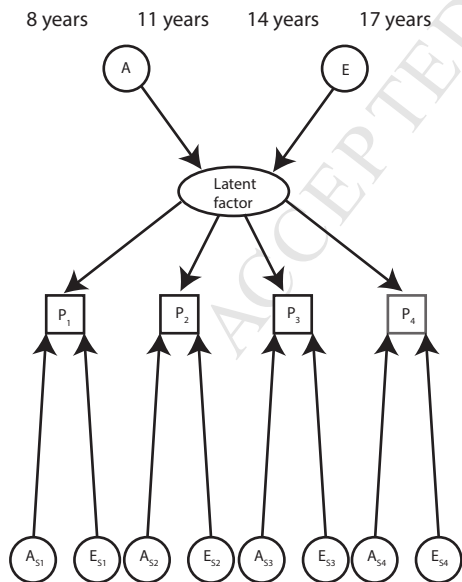
A



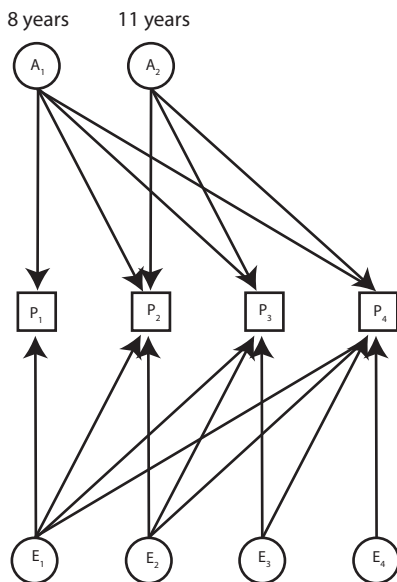
B



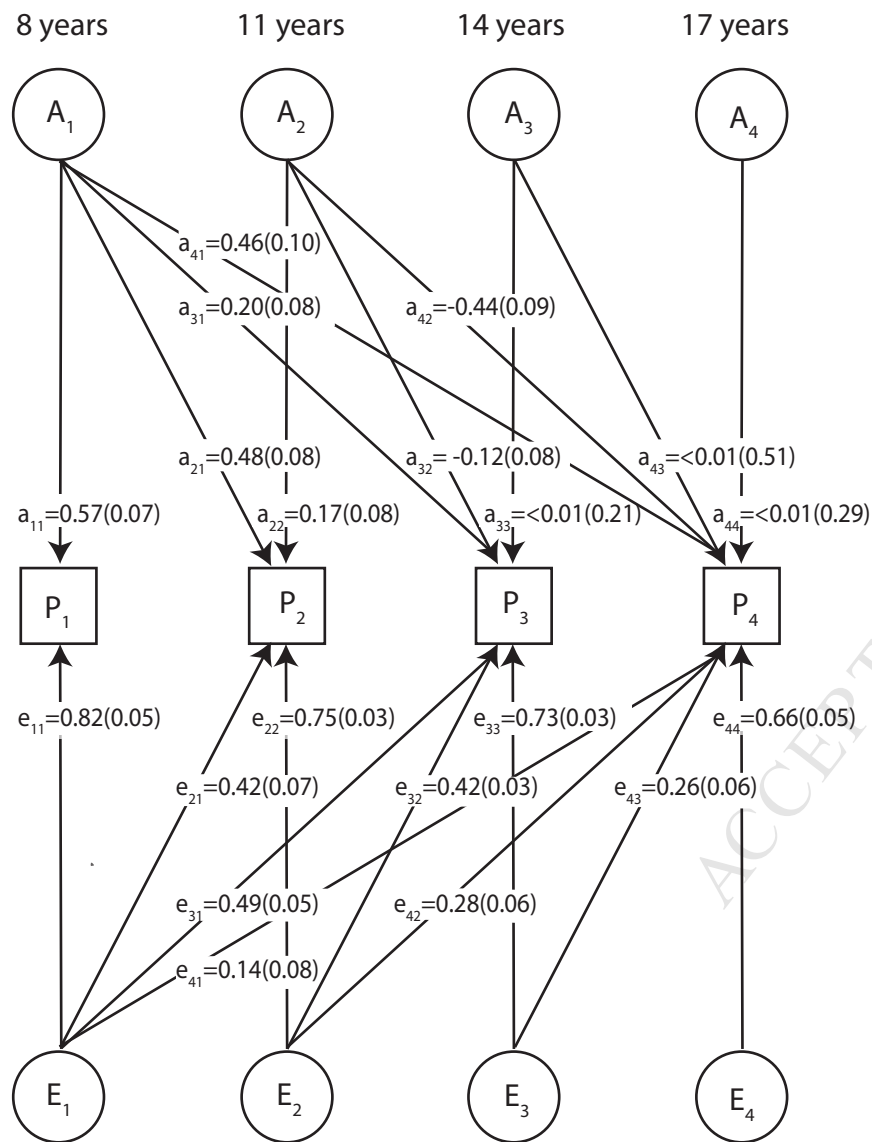
C



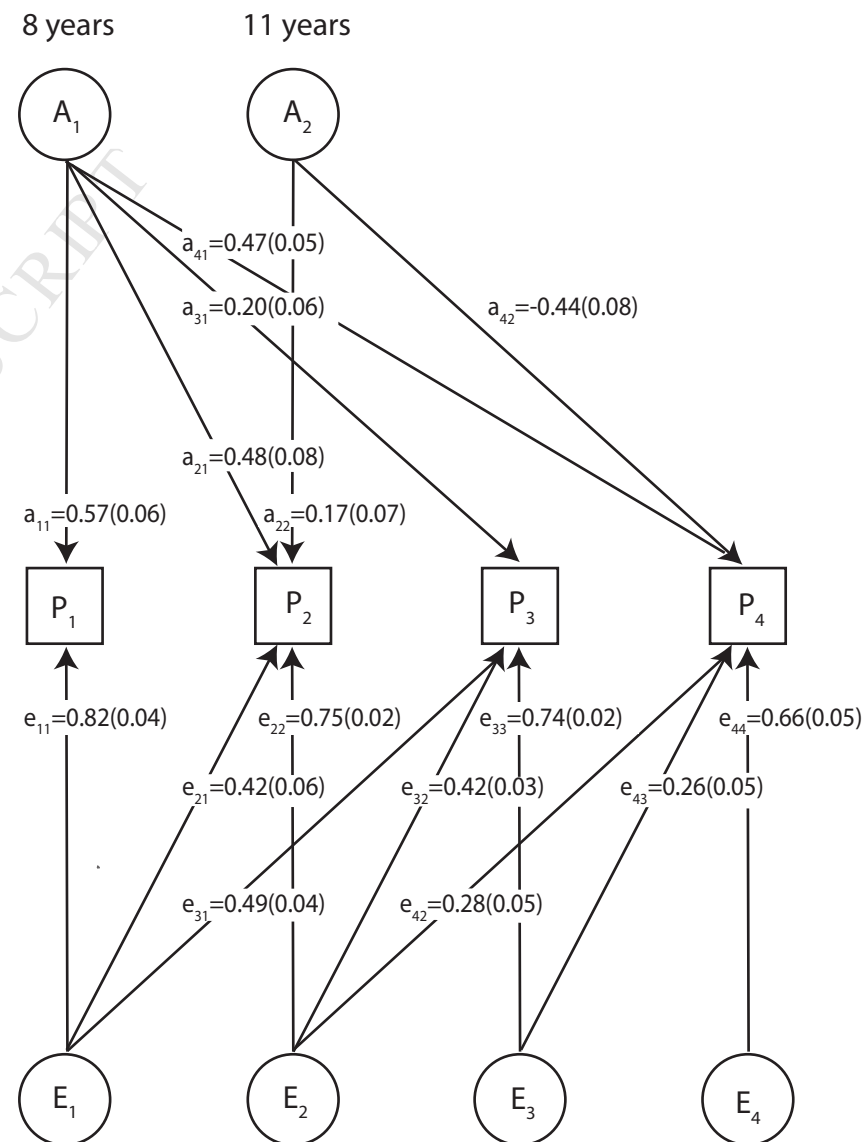
D



A



B



Developmental Changes Within the Genetic Architecture of Social Communication Behavior: A Multivariate Study of Genetic Variance in Unrelated Individuals

Supplementary Information

Supplementary Methods

- Genome-wide genotype information in the Avon Longitudinal Study of Parents and Children (ALSPAC)
- Genetic-relationship-matrix structural equation models (GSEM)
- Data simulation
- Attrition within ALSPAC
- Supplementary references
- Web resources
- R gsem package installation

Supplementary Tables

- Supplementary Table S1: Descriptives of SCDC scores
- Supplementary Table S2: Phenotypic correlation of SCDC scores
- Supplementary Table S3: Bivariate simulations
- Supplementary Table S4: Computational requirements (Bivariate simulations)
- Supplementary Table S5: Trivariate simulation
- Supplementary Table S6: Univariate GSEM of SCDC scores
- Supplementary Table S7: Univariate analysis of SCDC scores: GCTA(GREML) versus GSEM
- Supplementary Table S8: Multivariate GSEM of SCDC scores: Standardised factor loadings
- Supplementary Table S9: Multivariate GSEM of SCDC scores: Estimated genetic variances and bivariate correlations
- Supplementary Table S10: Bivariate analysis of SCDC scores: GCTA(GREML) versus GSEM
- Supplementary Table S11: Genetic correlation between SCDC scores and subsequent attrition
- Supplementary Table S12: Genetic correlations between SCDC attrition scores

Supplementary Figures

- Supplementary Figure S1: Path diagrams for simulated data sets
- Supplementary Figure S2: Bivariate simulation analyses
- Supplementary Figure S3: Bivariate genetic correlations between SCDC scores during development (bivariate GREML versus multivariate GSEM)

Supplementary MethodsGenome-wide genotype information in the Avon Longitudinal Study of Parents and Children (ALSPAC)

ALSPAC children were genotyped using the Illumina HumanHap550 quad chip genotyping platforms. The ALSPAC GWAS data was generated by Sample Logistics and Genotyping Facilities at the Wellcome Trust Sanger Institute and LabCorp (Laboratory Corporation of America) using support from 23andMe. After quality control (individual call rate > 0.97, SNP call rate > 0.95, minor allele frequency (MAF) > 0.01, Hardy-Weinberg equilibrium (HWE) $p > 10^{-7}$, and removal of individuals with cryptic relatedness and non-European ancestry), 8,237 children and 477,482 directly genotyped single nucleotide polymorphisms (SNPs) were kept within the study.

Genetic-relationship-matrix structural equation models (GSEM)

Similar to genetic restricted maximum likelihood (GREML) as implemented in genome-wide complex trait analysis (GCTA) software (1), GSEM use the genetic similarity between unrelated individuals to partition the expected phenotypic variance/covariance matrix into genetic and residual components. The model assumes that genetic and residual effects are independent, and that residual effects of different individuals are independent. A normally distributed phenotype in N unrelated individuals can thus be modelled (2) as

$$P \sim N(\mathbf{0}, \mathbf{G} \sigma_g^2 + \mathbf{I}(1 - \sigma_g^2)) \quad (1)$$

where P is a $N \times 1$ vector of phenotypes, \mathbf{G} is the $N \times N$ genetic correlation matrix of pairwise genome-wide genetic correlations between unrelated individuals, and \mathbf{I} is a $N \times N$ identity matrix. As in the GCTA software package (1), \mathbf{G} is the genetic relationship matrix (GRMs) constructed from common variants present on SNP chips, and σ_g^2 is an estimate of the genetic variance captured by these SNPs, while $(1 - \sigma_g^2)$ is an estimate of the residual variance (σ_e^2). Thus, similar to GREML (1), the total amount of phenotypic variance captured by genotyped SNPs can be estimated by fitting a univariate GSEM. GSEM uses full information maximum likelihood (FIML) and combinations of latent factor loadings and/or factor variances which can then be used to derive estimates of genetic and residual variances, covariances and correlations (see below) (3). More generally, the statistical framework of GSEM is analogous to twin analysis methodologies (4), where SEM (3) in genetically informative samples (with known average degree of genetic resemblance) are used to model the phenotypic covariance structure (5). In twin studies, genetic and environmental influences are parametrised as latent factors. The phenotypic covariance structure is often modelled by one or more additive genetic factors A (i.e. the total additive genetic effects), one or more common environmental factors C (i.e. environmental influences affecting the phenotype in

family members in an identical way) and one or more specific environmental factors E (i.e. unique exposure of family members to environmental factors). Instead of expected genetic correlations between twin pairs based on biometrical theory (4), GSEM uses genetic relationship matrices (GRMs) for genetic covariance structure modelling. Like GREML, it describes one or more additive genetic factors A and one or more residual factors E.

Within this study, we applied, univariate GSEM and multivariate GSEM, analogous to twin analysis (4). Assuming multivariate normality, and expressing the phenotype of each individual i as a deviation from the grand mean (5), the likelihood L_i for each person can be expressed as

$$\log(L_i) = -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2} P_i' \Sigma_i^{-1} P_i + c = -\frac{1}{2} \log |\Sigma_i| - \frac{1}{2} \text{tr}(P_i P_i' \Sigma_i^{-1}) \quad (2)$$

where Σ_i is the predicted variance/covariance matrix and P_i is the vector of phenotypes for the i^{th} individual with a grand mean of 0 ($P_i P_i'$ is the sample covariance matrix), and c is a constant term. The log likelihood (L) is then the sum of the log likelihoods for each individual.

$$\log(L) = \sum \log(L_i) \quad (3)$$

A saturated AE model can be obtained through a full decomposition of the genetic variance and residual variance into as many latent factors as there are observed measures (Cholesky decomposition). The Cholesky decomposition of the genetic variance can be described as follows (6): For a longitudinally assessed trait P with t repeat measurements, the first phenotypic measure, P_1 , is influenced by a latent genetic factor (A_1), which can also explain variance in the second and all following measures (P_2, \dots, P_t). The second measure (P_2) is, in addition, influenced by a second latent genetic factor A_2 , explaining phenotypic variance in P_2 and all following measures (P_3, \dots, P_t) not yet captured by A_1 , and so forth. The last measure (P_t) is, beside the latent genetic factors (A_1, \dots, A_{t-1}), influenced by a genetic

factor A_i , which does not explain variance within any of the previous measures (P_1, \dots, P_{i-1}) (4). We annotate the genetic factor loadings a (path coefficients) such that the first number indicates the direction of the effect (the variable to which the arrow points) and the second the origin of the effect (4).

The expected phenotypic covariance matrix for Z-standardised traits based on the factor model is

$$\Sigma = \lambda \Phi \lambda' + \Psi^2 \quad (4)$$

where λ is a lower triangular matrix of genetic factor loadings, Φ is a diagonal matrix of latent genetic factor variances (standardised to unit variance) such that Φ is an identity matrix I , and Ψ^2 a covariance matrix of residual influences (5). It is also possible to decompose the residual variance into latent residual factors, such that

$$\Sigma = \lambda \Phi \lambda' + \zeta \Theta \zeta' \quad (5)$$

where ζ is a lower triangular matrix of residual factor loadings and Θ is a diagonal matrix of latent residual factor variances (standardised to unit variance) such that Θ is an identity matrix I . For example, for a bivariate trait consisting of measures P_1 and P_2 , assuming two genetic factors (A_1 and A_2) and two genetic factors (E_1 and E_2), the expected phenotypic covariance matrix can be expressed as follows:

$$\Sigma = \begin{bmatrix} \sigma_{p_1}^2 & \sigma_{p_{12}} \\ \sigma_{p_{12}} & \sigma_{p_2}^2 \end{bmatrix} \quad (6)$$

with the relevant matrices

$$\lambda = \begin{bmatrix} a_{11} & 0 \\ a_{21} & a_{22} \end{bmatrix}, \Phi = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \zeta = \begin{bmatrix} e_{11} & 0 \\ e_{21} & e_{22} \end{bmatrix}, \Theta = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (7)$$

where $\sigma_{p_1}^2$ and $\sigma_{p_2}^2$ represent the phenotypic variances and $\sigma_{p_{12}}$ the phenotypic covariance.

The bivariate AE Cholesky decomposition of two standardised measures, as described above, can be visualised by means of a path diagram (Supplementary Figure S1A) and the expected phenotypic variances and covariances can be expressed as follows:

$$\sigma_{p_1}^2 = \sigma_{g_1}^2 + \sigma_{e_1}^2 = a_{11}^2 + e_{11}^2 = 1 \quad (8)$$

$$\sigma_{p_2}^2 = \sigma_{g_2}^2 + \sigma_{e_2}^2 = (a_{21}^2 + a_{22}^2) + (e_{21}^2 + e_{22}^2) = 1 \quad (9)$$

$$\sigma_{p_{12}} = \sigma_{g_{12}} + \sigma_{e_{12}} = a_{11}a_{21} + e_{11}e_{21} \quad (10)$$

where $\sigma_{g_1}^2$ and $\sigma_{g_2}^2$ represent the genetic variances and $\sigma_{g_{12}}$ the genetic covariance, and $\sigma_{e_1}^2$ and $\sigma_{e_2}^2$ the residual variances and $\sigma_{e_{12}}$ the residual covariance. The variance of the latent factors A_1 and A_2 , and E_1 and E_2 has been standardised to unit variance and is not shown.

Estimated genetic variances and covariances can subsequently be utilised to derive genetic correlations (GSEM- r_g) between two phenotypes (7), i.e. the extent to which two phenotypes share genetic factors (ranging from -1 to 1):

$$\rho_g = \frac{\sigma_{g_{12}}}{\sqrt{\sigma_{g_1}^2 \sigma_{g_2}^2}} \quad (11)$$

where $\sigma_{g_{12}}$ is the genetic covariance between phenotypes 1 and 2 and ρ_g the genetic correlation.

We fitted in this work an AE Cholesky decomposition model as baseline model and not, as commonly selected in twin research, a fully saturated model. A twin design, however, inherently contains genetic information based on phenotypic twin correlations in monozygotic versus dizygotic twins. A cohort sample consisting of unrelated individuals does not. Thus, fitting a fully saturated model in a general population sample of unrelated individuals will not provide information on genetic effects within that sample.

The goodness-of-fit of GSEM to empirical data was assessed using likelihood ratio test (LRT), the Akaike Information Criterion (AIC) (8) and the Bayesian Information Criterion (BIC) (9). The LRT is based on the difference in the negative log-likelihood ($-2LL$) of *a priori* defined models (model 2 and 3) and the saturated model (model 1), which is asymptotically chi-squared distributed with degrees of freedom equal to the difference in parameters between the models. AIC fit indices were calculated as

$$AIC = -2LL + 2k \quad (12)$$

where LL is the log-likelihood and k is the number of free model parameters in the model, with lower AIC values indexing a better model fit (8). The BIC indices take both goodness-of-fit and parsimony of the model into account, and lower BIC values indicate a better model fit (10). The index is defined as

$$BIC = -2LL + k \log (N) \quad (13)$$

where N is the number of independent observations.

GSEM were implemented within R (Rv3.2.4) via the `optim` function (stats library).

Data simulation

To evaluate the accuracy of multivariate GSEM, we carried out data simulations. Assuming multivariate normality, we simulated bivariate traits with two repeated measures (i.e. two genetic factors with their variances and their covariance; two residual factors with their variances and their covariance, Supplementary Figure S1A), assuming 5000 individuals and 20,000 SNPs per genetic factor, for 10 replicates. Phenotypic variances and covariances were estimated from genetic (a) and residual (e) factor loadings as expected under an AE Cholesky decomposition model (Supplementary Figure S1A). The simulated values are detailed in Supplementary Table S3. Across replicates, we calculated the mean average deviation (MAD; i.e. the absolute deviation from the mean), the root mean squared error (RMSE, i.e. the square root of the average squared difference between each replicate estimate and the true (simulated) value), and the squared bias (Bias², i.e. the squared difference between the mean of the replicate estimates and the true simulated value). As a benchmark test, we also carried out a trivariate trait simulation with three repeated measures (i.e. three genetic factors with their variances and their covariance and three residual factors with their variances and their covariance, Supplementary Figure S1B) assuming 5000 individuals and 20,000 SNPs per genetic factor, for one replicate.

Attrition within ALSPAC

To study non-participation within ALSPAC, analysis was restricted to participants who were alive at one year of age and had information on genome-wide data available (N=7,758). Dichotomic SCDC-missingness was defined as availability of mother-reported scores at 8, 11, 14 and 17 years of age. We estimated genetic correlations between these attrition scores as well as between SCDC scores and subsequent SCDC-missingness. This is possible as genetic correlations are independent of an underlying liability scale. For simplicity, we used bivariate GREML and not multivariate GSEM for the analysis, as there was little evidence for a genetic correlation between SCDC scores and subsequent attrition (see Results).

Supplementary references

1. Yang J, Lee SH, Goddard ME, Visscher PM (2011): GCTA: A Tool for Genome-wide Complex Trait Analysis. *Am. J. Hum. Genet.* 88: 76–82.
2. Golan D, Lander ES, Rosset S (2014): Measuring missing heritability: inferring the contribution of common variants. *Proc. Natl. Acad. Sci. U. S. A.* 111: E5272-5281.
3. Bollen KA (1989): *Structural Equations with Latent Variables*, 1 edition. New York: Wiley-Blackwell.
4. Neale M, Maes HHM (2004): *Methodology for genetic studies of twins and families.*, Dordrecht: Kluwer Academic Publishers.
5. Martin NG, Eaves LJ (1977): The genetical analysis of covariance structure. *Heredity* 38: 79–95.
6. Cherny SS (2005): Cholesky Decomposition. *Encycl. Stat. Behav. Sci.* John Wiley & Sons, Ltd.
7. Falconer PDS, Mackay PTF (1995): *Introduction to Quantitative Genetics*, 4 edition. Essex, England: Longman.
8. Akaike H (1987): Factor analysis and AIC. *Psychometrika* 52: 317–332.
9. Schwarz G (1978): Estimating the Dimension of a Model. *Ann. Stat.* 6: 461–464.
10. Maindonald J, Braun WJ (2010): *Data Analysis and Graphics Using R: An Example-Based Approach*, 3 edition. Cambridge ; New York: Cambridge University Press.
11. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, et al. (2010): Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* 42: 565–569.

Web resources

ALSPAC: <http://www.bris.ac.uk/alspac/researchers/data---access/data---dictionary>

HPC: <http://www.bristol.ac.uk/earthsciences/about/facilities/hpc.html>

PLINK2: <http://pngu.mgh.harvard.edu/~purcell/plink/plink2.shtml>

OPENMX: <http://openmx.psyc.virginia.edu/>

R: <https://cran.r-project.org/>

GCTA: <http://cnsgenomics.com/software/gcta/>

GSEM: <https://gitlab.gwdg.de/beate.stpourcain/gsem>

ACCEPTED MANUSCRIPT

R gsem package installation

```
install.packages("devtools")
devtools::install_github("hadley/devtools")
install.packages("msm")
install.packages("numDeriv")
devtools::install_git('https://gitlab.gwdg.de/beate.stpourcain/gsem')
```

#Note that ssl certificate issues during the installation can arise if the ca-certificates package on the client server is out of date

ACCEPTED MANUSCRIPT

Supplementary Tables

Supplementary Table S1: Descriptives of SCDC scores

Trait	Age(years)[range]	Male/Female	Mean(SD) [range]	Kurtosis	Skewness	N
8	7.7(0.14)[7.5;9.3]	2842/2709	2(3.71)[0;24]	9.12	2.19	5551
11	10.7(0.13)[10.5;13.8]	2751/2709	1(3.51)[0;24]	10.53	2.46	5460
14	13.9(0.15)[13.7;16.1]	2529/2531	1(3.59)[0;24]	9.08	2.20	5060
17	16.8(0.36)[16.5;18.3]	2024/2150	1(3.79)[0;24]	7.09	1.89	4174

SCDC - Social and Communication Disorders Checklist; The kurtosis for the standard normal distribution is 3 and the skewness is 0

Supplementary Table S2: Phenotypic correlation of SCDC scores

Age in years				
	8	11	14	17
8	1.00	0.61	0.50	0.38
11	0.57	1.00	0.56	0.41
14	0.49	0.57	1.00	0.51
17	0.39	0.45	0.56	1.00

SCDC - Social and Communication Disorders Checklist

Lower triangle: Spearman's rank correlation using pairwise complete observations; Upper triangle: Pearson product moment correlation using rank-transformed scores adjusted for age, sex and the two most significant ancestry-informative principal components

Supplementary Table S3: Bivariate simulations

Label	Sim	GCTA (GREML)			GSEM (FIML)			OpenMx 2.5 (FIML)		
		Mean(RMSE)	MAD	Bias ²	Mean(RMSE)	MAD	Bias ²	Mean(RMSE)	MAD	Bias ²
Var_{g1}	0.25	0.263(0.053)	0.039	1.8x10 ⁻⁴	0.263(0.053)	0.039	1.8x10 ⁻⁴	0.263(0.054)	0.039	1.8x10 ⁻⁴
Var_{g2}	0.50	0.48(0.091)	0.077	3.6x10 ⁻⁴	0.481(0.091)	0.077	3.5x10 ⁻⁴	0.481(0.091)	0.077	3.5x10 ⁻⁴
Cov_g	0.25	0.251(0.059)	0.055	4.6x10 ⁻⁷	0.251(0.059)	0.055	3.5x10 ⁻⁷	0.251(0.059)	0.055	3.0x10 ⁻⁷
Var_{e1}	0.75	0.740(0.05)	0.038	9.0x10 ⁻⁵	0.74(0.05)	0.038	9.4x10 ⁻⁵	0.74(0.05)	0.038	9.2x10 ⁻⁵
Var_{e2}	0.50	0.515(0.086)	0.071	1.8x10 ⁻⁴	0.513(0.086)	0.071	1.7x10 ⁻⁴	0.513(0.086)	0.071	1.7x10 ⁻⁴
Cov_e	0.087	0.085(0.052)	0.048	5.5x10 ⁻⁷	0.087(0.052)	0.048	7.0x10 ⁻⁷	0.087(0.052)	0.048	5.7x10 ⁻⁷
a₁₁	0.50	-	-	-	0.51(0.055)	0.041	1.1x10 ⁻⁴	0.51(0.055)	0.042	1.0x10 ⁻⁴
a₂₁	0.50	-	-	-	0.49(0.096)	0.082	9.7x10 ⁻⁵	0.49(0.096)	0.082	9.5x10 ⁻⁵
a₂₂	0.50	-	-	-	0.481(0.036)	0.025	3.8x10 ⁻⁴	0.481(0.036)	0.025	3.8x10 ⁻⁴
e₁₁	0.87	-	-	-	0.86(0.029)	0.021	3.6x10 ⁻⁵	0.86(0.029)	0.022	3.6x10 ⁻⁵
e₂₁	0.10	-	-	-	0.101(0.059)	0.054	8.3x10 ⁻⁷	0.101(0.059)	0.054	6.5x10 ⁻⁷
e₂₂	0.70	-	-	-	0.705(0.054)	0.044	2.1x10 ⁻⁵	0.704(0.054)	0.044	2.0x10 ⁻⁵

Bivariate trait data with two repeated measures (1 and 2, N=5000 for each trait) were simulated for 10 replicates (Supplementary Figure S1A): Absolute genetic factor loadings (a) are given with respect to two simulated genetic factors A₁ and A₂; Absolute residual factor loadings (e) are given with respect to two simulated residual factors E₁ and E₂; Simulated data were analysed with bivariate GCTA(GREML), bivariate GSEM(FIML) (full Cholesky factorisation) and bivariate OpenMx SEM (full Cholesky factorisation using the fastest OpenMx SEM algorithm, OpenMx(FIML) v2.5; see Supplementary Table S4).

Bias² - Squared bias defined as the squared difference between the mean of the replicate estimates and the true simulated value; Cov_e - Residual covariance; Cov_g - Genetic covariance; FIML - Full information maximum likelihood; GCTA - Genome-wide complex trait analysis; GREML - Genetic-relationship matrix restricted maximum likelihood; GSEM - Genetic-relationship-matrix structural equation models; MAD - Mean absolute deviation; Mean - Mean of 10 replicate estimates; RMSE - Root mean squared error; SEM - Structural equation models; Sim - Simulated (true) value; Var_e - Residual variance; Var_g - Genetic variance

Supplementary Table S4: Computational requirements (Bivariate simulations)

N	GCTA (GREML)		GSEM (FIML)		OpenMx 2.5 (FIML) ^a		OpenMx 2.7 (FIML) ^a		OpenMx 2.7 (mxGREML) ^a	
	RAM(GB)	Time(min)	RAM(GB)	Time(min)	RAM(GB)	Time(min)	RAM(GB)	Time(min)	RAM(GB)	Time(min)
1000	<0.1	<1	0.5	9	3.3	9	1.0	29	0.7	32
2000	1.0	2	1.5	98	12.8	44	3.7	213	2.4	298
3000	2.1	9	3.2	121	28.4	115	8.3	600	5.2	785
4000	4.0	25	5.6	298	50.2	204	16.0	1260	9.6	1681
5000	5.9	28	8.7	301	78.4	351	22.8	2694	15.1	2597

Bivariate trait data with two repeated measures were simulated for one replicate using a series of different sample sizes (1000 to 5000). The estimated parameter values are identical to those described in Supplementary Table S3. Analyses were conducted using one core (2.60 GHz) with access to 64 GB or, if required, 256 GB; RAM(GB) - Memory in Giga Bytes; Time(min) - Time in minutes; a - All OpenMx (mx) analyses were carried out with the options: options(mxCondenseMatrixSlots=TRUE) and mxOption(NULL,"Default optimizer","NPSOL"); GCTA - Genome-wide complex trait analysis; GREML - Genetic-relationship matrix restricted maximum likelihood, FIML - Full information maximum likelihood

Supplementary Table S5: Trivariate simulation

Label	Sim	GSEM (FIML) Estimate(SE)	OpenMx 2.5 (FIML) Estimate(SE)
a11	0.80	0.83(0.043)	0.83(0.043)
a21	0.50	0.498(0.055)	0.498(0.055)
a31	0.00	-0.010(0.062)	-0.011(0.062)
a22	0.40	0.403(0.04)	0.403(0.04)
a32	0.50	0.465(0.071)	0.464(0.071)
a33	0.10	0.101(0.165)	0.098(0.17)
e11	0.60	0.572(0.058)	0.572(0.058)
e21	0.60	0.625(0.06)	0.625(0.06)
e31	0.20	0.211(0.086)	0.212(0.086)
e22	0.48	0.458(0.04)	0.458(0.04)
e32	0.50	0.537(0.067)	0.538(0.067)
e33	0.67	0.682(0.026)	0.682(0.026)

A trivariate trait with three standardised measures (1, 2 and 3) was simulated (a single replicate) as a benchmark test (Supplementary Figure S1B): Absolute genetic factor loadings (a) are given with respect to three simulated genetic factors A₁, A₂ and A₃; Absolute residual factor loadings (e) are given with respect to three simulated residual factors E₁, E₂ and E₃; Simulated data were analysed with trivariate GSEM(FIML) (full Cholesky factorisation) and trivariate OpenMx SEM (full Cholesky factorisation using the fastest OpenMx SEM algorithm, OpenMx(FIML) v2.5; see Supplementary Table S4).

SEM - Structural equation models; Sim - Simulated (true) value; GSEM - Genetic-relationship-matrix structural equation models

Supplementary Table S6: Univariate GSEM of SCDC scores

Age(y)	Var _g (SE)	GSEM	
		<i>p</i>	N
8	0.25(0.06)	3.36x10 ⁻⁵	5551
11	0.22(0.06)	2.94x10 ⁻⁴	5460
14	0.086(0.06)	0.18	5060
17	0.47(0.09)	4.40x10 ⁻⁸	4174

GSEM - Genetic-relationship-matrix structural equation models; Var_g - Additive genetic variance; SCDC - Social and Communication Disorders Checklist

Univariate GSEM was carried out by estimating the genetic variance as a variance component without constraints (2-sided test). The reported Var_g estimates are equivalent to SNP-h² estimates due to the standardisation of the analysed traits.

Supplementary Table S7: Univariate analysis of SCDC scores: GCTA(GREML) versus GSEM

Age	GCTA(GREML)			GSEM		
	Var _g (SE)	<i>p</i>	N	Var _g (SE)	<i>p</i>	N
8	0.23(0.07)	1.6x10 ⁻⁴	4971	0.23(0.07)	1.6x10 ⁻⁴	4971
11	0.15(0.07)	0.011	4895	0.15(0.07)	0.011	4895
14	0.097(0.07)	0.077	4566	0.099(0.07)	0.077	4566
17	0.47(0.09)	5.0x10 ⁻⁹	3779	0.47(0.09)	5.0x10 ⁻⁹	3779

Age - Age in years; GCTA - Genome-wide complex trait analysis; GREML - Genetic-relationship matrix restricted maximum likelihood; GSEM - Genetic-relationship-matrix structural equation models; SCDC - Social and Communication Disorders Checklist; Var_g - Additive genetic variance

Differences compared with the total sample N are due to the exclusion of individuals with a relatedness of $\geq 2.5\%$ from all ALSPAC participants with genome-wide data, allowing for identical sample numbers across GREML and GSEM analyses. The reported Var_g estimates are equivalent to SNP- h^2 estimates due to the standardisation of the analysed traits.

Univariate GSEM was carried out by estimating the genetic variance without constraints. Note, that the GREML likelihood ratio test, as implemented in GCTA (11), follows by default a 50:50 mixture distribution with a point mass at 0 and a chi-squared distribution (df=1), which is comparable to a one-tailed *p*-value. For comparison with GSEM, however, an unconstrained GCTA option was selected.

Supplementary Table S8: Multivariate GSEM of SCDC scores: Standardised factor loadings

Label	Full Cholesky decomposition (Model 1)		Best-fitting model (Model 5)	
	Estimate(SE)	<i>p</i>	Estimate (SE)	<i>p</i>
a11	0.57(0.07)	5.09x10 ⁻¹³	0.57(0.06)	8.69x10 ⁻¹⁷
a21	0.48(0.08)	1.30x10 ⁻⁸	0.48(0.08)	2.57x10 ⁻⁹
a31	0.20(0.08)	2.15x10 ⁻²	0.20(0.06)	4.10x10 ⁻³
a41	0.46(0.10)	2.15x10 ⁻⁶	0.47(0.05)	3.29x10 ⁻²¹
a22	0.17(0.08)	4.09x10 ⁻²	0.17(0.07)	1.94x10 ⁻²
a32	-0.12(0.08)	1.69x10 ⁻¹	-	-
a42	-0.44(0.09)	2.60x10 ⁻⁶	-0.44(0.08)	8.24x10 ⁻⁸
a33	<0.01(0.21)	1	-	-
a43	<0.01 (0.51)	1	-	-
a44	<0.01 (0.29)	1	-	-
e11	0.82(0.05)	<10 ⁻¹⁰	0.82(0.04)	<10 ⁻¹⁰
e21	0.42(0.07)	1.73x10 ⁻⁹	0.42(0.06)	<10 ⁻¹⁰
e31	0.49(0.05)	<10 ⁻¹⁰	0.49(0.04)	<10 ⁻¹⁰
e41	0.14(0.08)	6.69x10 ⁻²	-	-
e22	0.75(0.03)	<10 ⁻¹⁰	0.75(0.02)	<10 ⁻¹⁰
e32	0.42(0.03)	<10 ⁻¹⁰	0.42(0.03)	<10 ⁻¹⁰
e42	0.28(0.06)	1.13x10 ⁻⁶	0.28(0.05)	1.55x10 ⁻⁷
e33	0.73(0.03)	<10 ⁻¹⁰	0.74(0.02)	<10 ⁻¹⁰
e43	0.26(0.06)	1.51x10 ⁻⁵	0.26(0.05)	8.28x10 ⁻⁹
e44	0.66(0.05)	<10 ⁻¹⁰	0.66(0.05)	<10 ⁻¹⁰

The full Cholesky decomposition model and its best-fitting reduced form are described in Table 1 (Model 1 and model 5 respectively) and Figure 3: Genetic factor loadings *a* are given with respect to the latent genetic factor A₁(8 years), factor A₂(11 years), factor A₃(14 years) and factor A₄(17 years) and residual factor loadings *e* with respect to factor E₁(8 years), factor E₂(11 years), factor E₃(14 years) and factor E₄(17 years). 3,295 participants had non-missing scores across all ages.

GSEM - Genetic-relationship-matrix structural equation models; SCDC - Social and Communication Disorders Checklist

Supplementary Table S9: Multivariate GSEM of SCDC scores: Estimated genetic variances and bivariate correlations

Full Cholesky decomposition (Model 1)				
Var_g(SE)				
Age	8	11	14	17
	0.32(0.08)	0.26(0.08)	0.05(0.03)	0.41(0.09)
r_g(SE)				
Age	8	11	14	17
8	1	-	-	-
11	0.95(0.05)	1	-	-
14	0.86(0.20)	0.65(0.29)	1	-
17	0.73(0.12)	0.46(0.16)	0.98(0.08)	1
Best-fitting model (Model 5)				
Var_g(SE)				
Age	8	11	14	17
	0.32(0.07)	0.26(0.07)	0.04(0.03)	0.41(0.07)
r_g(SE)				
Age	8	11	14	17
8	1	-	-	-
11	0.95(0.05)	1	-	-
14	1.00(<0.01)	0.95(0.05)	1	-
17	0.73(0.08)	0.46(0.13)	0.73(0.08)	1

The full Cholesky decomposition model and its best-fitting reduced form are described in Table 1 (Model 1 and model 5 respectively) and Figure 3. 3,295 participants had non-missing scores across all ages.

Age - Age at measurement in years; GSEM - Genetic-relationship-matrix structural equation models; r_g - Genetic correlation; SCDC - Social and Communication Disorders Checklist; Var_g - Genetic variance

Note that SCDC scores at 14 years were retained within the model, irrespective of their low SNP-heritability, due to their genetic correlations with other SCDC measures earlier and later during development. The reported Var_g estimates are equivalent to SNP-h² estimates and based on standardised path coefficients.

Supplementary Table S10: Bivariate analysis of SCDC scores: GCTA(GREML) versus GSEM

Full	Age	8 vs 11 Obs=9,866	8 vs 14 Obs =9,537	8 vs 17 Obs =8,750	11 vs 14 Obs =9,461	11 vs 17 Obs =8,674	14 vs 17 Obs =8,345
Method		Estimate(SE)	Estimate(SE)	Estimate(SE)	Estimate(SE)	Estimate(SE)	Estimate(SE)
GCTA(GREML)	A ₁	0.22(0.065)	0.22(0.067)	0.23(0.068)	0.17(0.067)	0.16(0.068)	0.09(0.07)
	A ₂	0.15(0.065)	0.09(0.069)	0.47(0.087)	0.11(0.069)	0.48(0.086)	0.49(0.086)
	Cov _g	0.17(0.053)	0.1(0.052)	0.15(0.056)	0.11(0.055)	0.11(0.057)	0.17(0.061)
	r _g	0.95(0.156)	0.68(0.266)	0.46(0.145)	0.79(0.217)	0.40(0.17)	0.82(0.257)
	<i>p</i> _{one-tailed}	0.00030	0.025	0.0025	0.018	0.024	0.00014
Non-missing	Age	8 vs 11 Obs =8,434	8 vs 14 Obs =7,872	8 vs 17 Obs =6,696	11 vs 14 Obs =8,820	11 vs 17 Obs =6,942	14 vs 17 Obs =7,014
Method		Estimate(SE)	Estimate(SE)	Estimate(SE)	Estimate(SE)	Estimate(SE)	Estimate(SE)
GCTA(GREML)	A ₁	0.16(0.075)	0.3(0.082)	0.22(0.092)	0.13(0.076)	0.11(0.088)	0.03(0.089)
	A ₂	0.08(0.077)	0.07(0.081)	0.41(0.099)	0.01(0.076)	0.47(0.096)	0.44(0.091)
	Cov _g	0.11(0.062)	0.12(0.064)	0.15(0.072)	0.04(0.062)	0.08(0.07)	0.11(0.071)
	r _g	1(0.354)	0.89(0.418)	0.51(0.188)	-	0.33(0.243)	-
	<i>p</i> _{one-tailed}	0.5	0.019	0.013	-	0.13	-
GSEM	A ₁	0.17(0.075)	0.3(0.081)	0.22(0.091)	0.13(0.074)	0.11(0.084)	0.04(0.037)
	A ₂	0.1(0.061)	0.07(0.078)	0.41(0.094)	0.02(0.054)	0.47(0.092)	0.47(0.092)
	Cov _g	0.13(0.06)	0.13(0.062)	0.15(0.069)	0.05(0.058)	0.08(0.067)	0.14(0.065)
	r _g	1(0)	0.89(0.422)	0.51(0.182)	-	0.33(0.229)	-

Age - Age at measurement in years; GCTA - Genome-wide complex trait analysis; GREML - Genetic-relationship matrix restricted maximum likelihood; GSEM - Genetic-relationship-matrix structural equation models; SCDC - Social and Communication Disorders Checklist

Bivariate analyses were conducted using bivariate GCTA(GREML) and bivariate GSEM (full Cholesky decomposition of the genetic and residual variance). Analyses were either carried out using the full phenotypic information (GCTA(GREML) only; including individuals with information at one time point only), or restricting the analysis to individuals with non-missing information across both time points. Individuals with a relatedness of $\geq 2.5\%$ were excluded. For comparisons, genetic variances (Var_g), covariances (Cov_g) and correlations (r_g) are shown. r_g estimates for traits with Var_g<0.05 are not reported due to large estimation errors.

Supplementary Table S11: Genetic correlation between SCDC scores and subsequent attrition

$r_g(\text{SE}), p_{\text{one-tailed}}$	8 (Score)	11 (Score)	14 (Score)	17 (Score)
8 (miss)	-	-	-	-
11 (miss)	0.33(0.23), $p=0.06$	-	-	-
14 (miss)	0.39(0.19), $p=0.02$	0.27(0.27), $p=0.15$	-	-
17 (miss)	0.24(0.20), $p=0.11$	0.20(0.28), $p=0.24$	0.28(0.34), $p=0.19$	-

miss - Missingness in Social and Communication Disorders Checklist (SCDC) scores; Scores - SCDC scores

Genetic correlations (r_g) were estimated using bivariate genetic-relationship matrix restricted maximum likelihood (GREML) as implemented in genome-wide complex trait analysis (GCTA) software

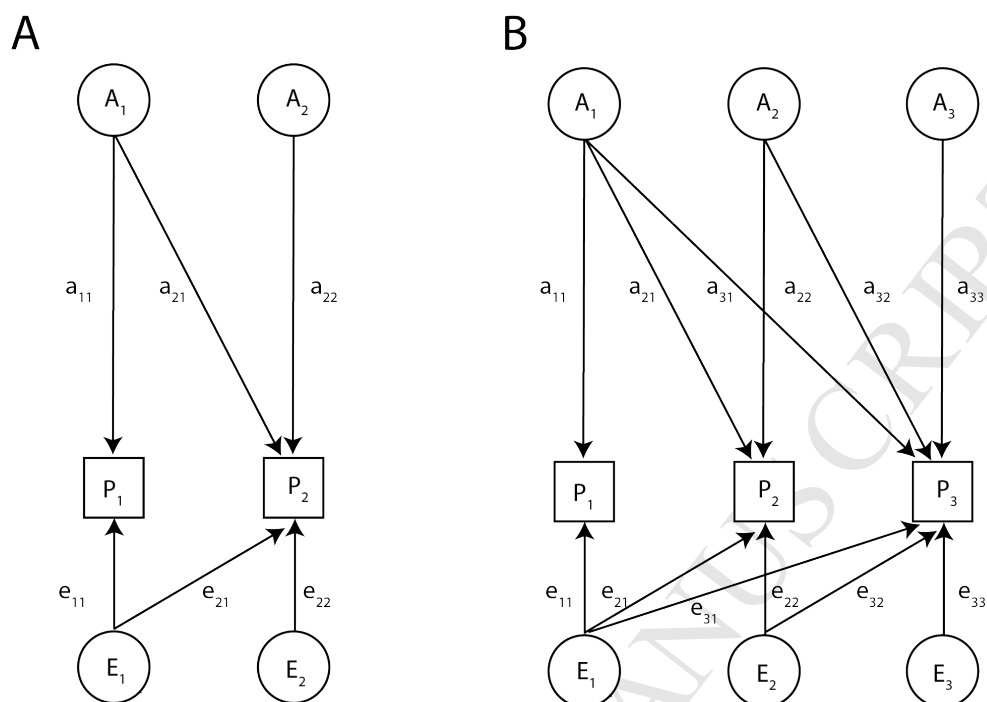
Supplementary Table S12: Genetic correlations between SCDC attrition scores

$r_g(\text{SE}), p_{\text{one-tailed}}$				
Miss	8	11	14	17
8	-	-	-	-
11	0.81(0.13), $p=4.8 \times 10^{-4}$	-	-	-
14	0.83(0.12), $p=2.2 \times 10^{-5}$	0.84(0.11), $p=1.1 \times 10^{-4}$	-	-
17	1.00(0.15), $p=8.4 \times 10^{-7}$	1.00(0.15), $p=2.2 \times 10^{-5}$	0.91(0.09), $p=3.3 \times 10^{-6}$	-

Miss - Missingness in Social and Communication Disorders Checklist (SCDC) scores

Genetic correlations (r_g) were estimated using bivariate genetic-relationship matrix restricted maximum likelihood (GREML) as implemented in genome-wide complex trait analysis (GCTA) software

Supplemental Figures

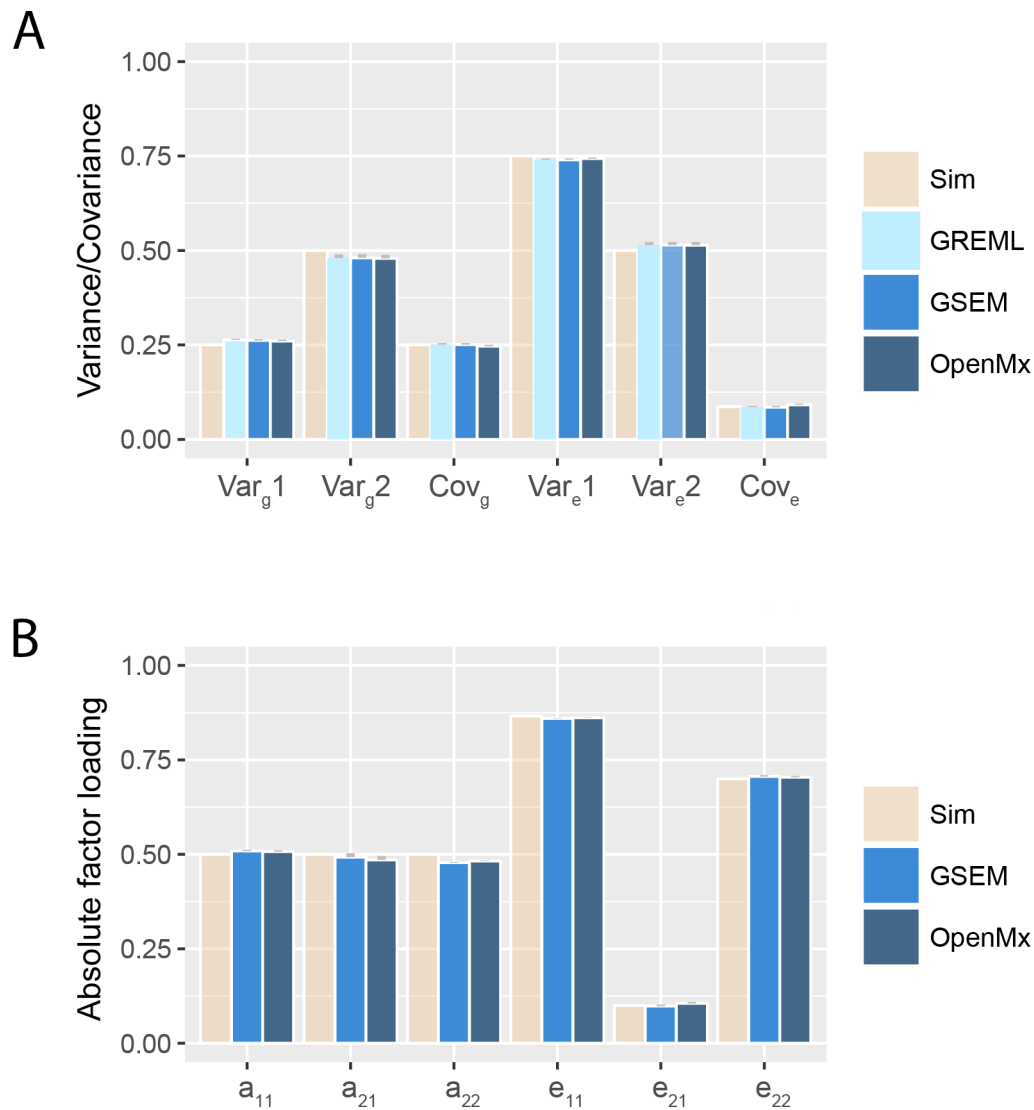


Supplementary Figure S1: Path diagrams for simulated data sets

A - A bivariate trait consisting of two standardised measures P_1 and P_2 was simulated for two genetic factors (A_1 and A_2) and two residual factors (E_1 and E_2), shown with genetic and residual factor loadings, assuming 5000 individuals and 20,000 SNPs per genetic factor (10 replicates)

B - A trivariate trait consisting of three standardised measures P_1 , P_2 and P_3 was simulated for three genetic factors (A_1 , A_2 and A_3) and three residual factors (E_1 , E_2 and E_3), shown with genetic and residual factor loadings, assuming 5000 individuals and 20,000 SNPs per genetic factor (one replicate)

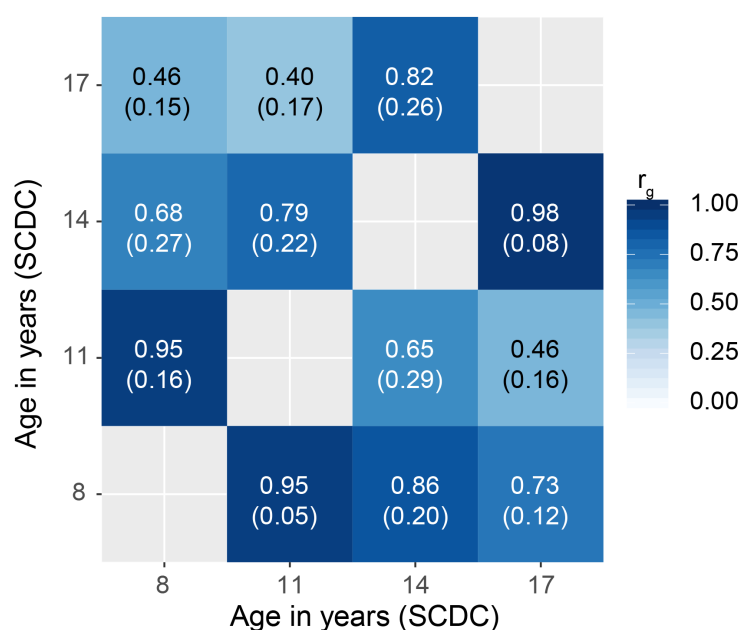
Observed phenotypic measures are represented by squares, while latent factors are represented by a circle. Single headed arrows ('paths') denote causal relationships between variables and are shown for genetic factor loadings (a) and residual factor loadings (e). Note that the variance of latent variables is constrained to unit variance, this is omitted from the diagrams to improve clarity.



Supplementary Figure S2: Bivariate simulation analyses

Mean estimated genetic and environmental variances and covariances (A) and factor loadings (B) are shown for 10 replicates with mean squared errors as grey bars. Simulated data were analysed with GCTA(GREML), GSEM (full Cholesky factorisation) and OpenMx SEM (full Cholesky factorisation, OpenMx 2.5, FIML). Bivariate traits with two standardised measures (1 and 2) were simulated assuming two genetic factors (A_1 and A_2) and two residual factors (E_1 and E_2 ; Supplementary Figure S1A, Supplementary Table S3).

a - Genetic factor loading; Cov_e - Residual covariance; Cov_g - Genetic covariance; e - Residual factor loading; FIML - Full information maximum likelihood; GCTA - Genome-wide complex trait analysis; GREML - Genetic-relationship-matrix residual maximum likelihood; GSEM - Genetic-relationship-matrix structural equation models; SEM - Structural equation models; Sim - Simulated value; Var_e - Residual variance; Var_g - Genetic variance



Supplementary Figure S3: Bivariate genetic correlations between SCDC scores during development (bivariate GREML versus multivariate GSEM)

Genetic correlations (r_g) were estimated with multivariate GSEM (lower triangle; as shown for the full Cholesky decomposition model in Supplementary Table S9, $N=13,180$ observations) and bivariate GCTA(GREML) (upper triangle; Supplementary Table S10, $N \leq 9,866$ observations; analyses using full phenotypic information) and are shown with their standard errors.

Note that both bivariate GREML and bivariate GSEM analyses using identical sample numbers provided nearly identical r_g estimates (Supplementary Table S10, non-missing sample analyses).

GCTA - Genome-wide complex trait analysis; GREML - Genetic-relationship matrix restricted maximum likelihood; GSEM - Genetic-relationship-matrix structural equation models