### Check for updates

### OPEN ACCESS

EDITED BY Roseann E. Peterson, Suny Downstate Health Sciences University, United States

REVIEWED BY Felix Christian Tropf, ENSAE ParisTech Ecole Nationale de la Statistique et de l'Administration Economique, France Qingwen Li, Chinese Academy of Sciences, China S. Mason Garrison, Wake Forest University, United States

\*CORRESPONDENCE Gang Chen, ⊠ gangchen@mail.nih.gov

RECEIVED 04 November 2024 ACCEPTED 06 March 2025 PUBLISHED 02 April 2025

#### CITATION

Chen G, Moraczewski D and Taylor PA (2025) Improving accuracy and precision of heritability estimation in twin studies through hierarchical modeling: reassessing the measurement error assumption. *Front. Genet.* 16:1522729. doi: 10.3389/fgene.2025.1522729

#### COPYRIGHT

© 2025 Chen, Moraczewski and Taylor. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

# Improving accuracy and precision of heritability estimation in twin studies through hierarchical modeling: reassessing the measurement error assumption

### Gang Chen<sup>1\*</sup>, Dustin Moraczewski<sup>2</sup> and Paul A. Taylor<sup>1</sup>

<sup>1</sup>Scientific and Statistical Computing Core, National Institute of Mental Health, Bethesda, MD, United States, <sup>2</sup>Data Science and Sharing Team, National Institute of Mental Health, Bethesda, MD, United States

**Introduction:** The conventional approach to estimating heritability in twin studies implicitly assumes either the absence of measurement error or that any measurement error is incorporated into the nonshared environment component. However, this assumption can be problematic when it does not hold or when measurement error cannot be reasonably classified as part of the nonshared environment.

**Methods:** In this study, we demonstrate the need for improvement in the conventional structural equation modeling (SEM) used for estimating heritability when applied to trait data with measurement errors. The critical issue revolves around an assumption concerning measurement errors in twin studies. In cases where traits are measured using samples, data is aggregated during preprocessing, with only a centrality measure (e.g., mean) being used for modeling. Additionally, measurement errors resulting from sampling are assumed to be part of the nonshared environment and are thus overlooked in heritability estimation. Consequently, the presence of intra-individual variability remains concealed. Moreover, recommended sample sizes are typically based on the assumption of no measurement errors.

**Results:** We argue that measurement errors in the form of intra-individual variability are an intrinsic limitation of finite sampling and should not be considered as part of the nonshared environment. Previous studies have shown that the intra-individual variability of psychometric effects is significantly larger than the inter-individual counterpart. Here, to demonstrate the appropriateness and advantages of our hierarchical linear modeling approach in heritability estimation, we utilize simulations as well as a real dataset from the ABCD (Adolescent Brain Cognitive Development) study. Moreover, we showcase the following analytical insights for data containing non-negligible measurement errors: i) The conventional SEM may underestimate heritability. ii) A hierarchical model provides a more accurate assessment of heritability. iii) Large samples, exceeding 100 observations or thousands of twins, may be necessary to reduce imprecision.

**Discussion:** Our study highlights the impact of measurement error on heritability estimation and introduces a hierarchical model as a more accurate alternative. These findings have significant implications for understanding individual differences and improving the design and analysis of twin studies.

### KEYWORDS

heritability, twin studies, ACE model, Falconer's method, intra-individual variability, hierarchical modeling, data generating mechanism, Bayesian statistics

## 1 Introduction

As an indication of potential predictability, heritability is an important concept in assessing individual differences. As the proportion of trait variability ascribed to genetics, heritability offers a unique perspective for quantifying the role of genetics in complex traits (Downes and Turkheimer, 2022; Robette et al., 2022). Twins provide a hypothetically well-controlled scenario where genetics, environment, and their interaction can be statistically separated and apportioned.

# 1.1 Heritability estimation: ACE model and Falconer's formula

Conventional twin studies are typically conceptualized with three hierarchies of data structure: individual, family, and zygosity. The individual measures are nested within families, which are further categorized as either monozygotic (MZ) or dizygotic (DZ) twins. A model can be formulated for a quantitative trait of interest that is measured at the individual level. In the popular ACE formulation (Maes, 2005; Downes and Matthews, 2020; Hunter, 2021), the trait data  $y_{i(f(z))}$  is expressed as the combination of latent components through the three indices of individual (i = 1, 2, ..., I), family (f = 1, 2, ..., F), and zygosity (z = MZ, DZ):

individual : 
$$y_i(f(z)) = \alpha + A_i(f(z)) + C_i(f(z)) + E_i(f(z))$$
. (1)

Each pair of parentheses indicates a nesting structure among the subscripts. The intercept  $\alpha$  captures the overall trait effect at the population level. The acronym for the ACE model reflects the three latent sources of variability.  $A_{i(f(z))}$  represents the additive genetic effects, and  $C_{i(f(z))}$  represents the common or shared environmental effects. In addition,  $E_{i(f(z))}$  characterizes the unique or nonshared environmental effects.

The variances associated with the three latent components are crucial parameters in twin studies. One may make the following assumptions for two twins  $i_1$  and  $i_2$  within a family f of zygosity z (Arbet et al., 2020),

$$\begin{bmatrix} A_{i_1} \\ A_{i_2} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho_z \\ \rho_z & 1 \end{bmatrix} \sigma_{A_z}^2\right), \begin{bmatrix} C_{i_1} \\ C_{i_2} \end{bmatrix}$$
$$\sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \sigma_{C_z}^2\right), \begin{bmatrix} E_{i_1} \\ E_{i_2} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \sigma_{E_z}^2\right).$$
(2)

The relatedness  $\rho_z$  for the additive genetic effects between two twins  $i_1$  and  $i_2$  in a family is 1 when z = MZ and 0 when z = DZ. As a side note, the main notations used in this paper are listed in Table 1.

A crucial feature of the Model 2 is the inclusion of the *homogeneity assumption*. This assumption is necessary when estimating variances within the model framework in (Model 2). Six variances, namely  $\sigma_{A_z}^2$ ,  $\sigma_{C_z}^2$ , and  $\sigma_{E_z}^2$  for z = DZ and MZ, need to be estimated, resulting in an undetermined system. To resolve this identifiability issue, the variances are assumed to be homogeneous across MZ and DZ twins (Arbet et al., 2020):

$$\sigma_{A_{\rm MZ}}^2 = \sigma_{A_{\rm DZ}}^2 = \sigma_A^2, \quad \sigma_{C_{\rm MZ}}^2 = \sigma_{C_{\rm DZ}}^2 = \sigma_C^2, \quad \sigma_{E_{\rm MZ}}^2 = \sigma_{E_{\rm DZ}}^2 = \sigma_E^2.$$
(3)

The assumption effectively reduces the number of variance parameters by half. However, this also leads to having homogeneity of total variance between MZ and DZ:  $\sigma_{A_{\text{MZ}}}^2 + \sigma_{C_{\text{MZ}}}^2 + \sigma_{E_{\text{MZ}}}^2 + \sigma_{E_{\text{DZ}}}^2 = \sigma_A^2 + \sigma_C^2 + \sigma_E^2 + \sigma_E^2$ .

The heritability in a twin study can be defined under the homogeneity assumption (Equation 3). The classical methodology is to adopt structural equation modeling (SEM) (e.g., Rijsdijk and Sham, 2002; Holst et al., 2016; Neale et al., 2016; Bates et al., 2019) to estimate the three variances. Then, the three proportions of total variability attributed to additive genetic effects, common environment, and nonshared environment are expressed respectively as,

$$h^{2} = \frac{\sigma_{A}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2}}, c^{2} = \frac{\sigma_{C}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2}}, e^{2} = \frac{\sigma_{E}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2}}.$$
 (4)

Effect decomposition under the SEM (Model 2) and the associated estimation of heritability can be visually represented as a path diagram (Figure 1), which is analogous to a directed acyclic diagram in causal inference. The values of  $h^2$ ,  $c^2$ , and  $e^2$  in common practice are typically reported with their point estimates, accompanied by their uncertainty expressed through standard errors or 95% uncertainty intervals<sup>1</sup>.

The three variability proportions of  $h^2$ ,  $c^2$ , and  $e^2$  can alternatively be expressed as the relatedness between the twins of each zygosity. We denote  $r_{MZ}$  and  $r_{DZ}$  as the correlations between two twins  $i_1$  and  $i_2$  within a family f for MZ and DZ, respectively. The following can be derived per the ACE model under the SEM Formulation 2,

$$r_{z} = \operatorname{corr}\left(y_{i_{1}}(f(z)), y_{i_{2}}(f(z))\right) = \frac{\rho_{z}\sigma_{A}^{2} + \sigma_{C}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2}} = \rho_{z}h^{2} + c^{2}, z = MZ, DZ.$$
(5)

The expression (Equation 5) leads to Falconer's formula (Falconer and MacKay, 1996):

$$h^2 = 2(r_{\rm MZ} - r_{\rm DZ}), \quad c^2 = 2r_{\rm DZ} - r_{\rm MZ}, \quad e^2 = 1 - h^2 - c^2 = 1 - r_{\rm MZ}.$$
(6)

## 1.2 Motivation: addressing the limitations of the conventional SEM framework

We address the challenges of heritability estimation from the perspectives of *accuracy* and *precision*, which are typically the objectives of studies to maximize. Here, we define the accuracy of an estimate as the absence of systematic biases. Inaccuracy, for example, implies an upward or downward shift in a point estimate of effect centrality (e.g., mean, mode, median), an uncertainty interval, or a posterior distribution. On the other hand, we define imprecision as the extent of uncertainty in an estimate, which can be quantified as the standard error or an uncertainty interval.

An implicit assumption is present within the conventional SEM regarding measurement errors, the random or stochastic

<sup>1</sup> The term "uncertainty interval" is used here to include confidence intervals in frequentist usage and credible intervals in Bayesian methods.

TABLE 1 A reference table of major variables and parameters used with heritability modeling. Quantities which originate in hierarchical linear modeling (HLM) and structural equation modeling (SEM) are noted explicitly.

Term	Description
$\sigma^2_{A_z}$	(SEM) variance of the additive genetic effects, for a given zygosity $\boldsymbol{z}$
$\sigma^2_{C_z}$	(SEM) variance of the common (shared) environment effects, for a given zygosity $\boldsymbol{z}$
$\sigma_{E_z}^2$	(SEM) variance of the nonshared genetic effects, for a given zygosity $\boldsymbol{z}$
$\sigma_A^2, \sigma_C^2, \sigma_E^2$	(SEM) variance components (above) under homogeneity assumption (equal variance across zygosity)
$\rho_z$	(SEM) relatedness for additive genetic effects between two twins
h <sup>2</sup>	Heritability, proportion of total variability attributed to additive genetic effects
<i>c</i> <sup>2</sup>	Proportion of total variability attributed to common environment
<i>e</i> <sup>2</sup>	Proportion of total variability attributed to nonshared environment
rz	Correlation between two twins of zygosity $z$
$\sigma_z^2$	(HLM) individual-level trait variance for zygosity $\boldsymbol{z}$
$\tau_z^2$	(HLM) family-level variance for zygosity $z$
$\omega_0^2$	(HLM) theoretical cross-trial variance of observations; intra-individual variability
$\omega^2$	(HLM) cross-trial sampling variance
$R_{\rm v}^2$	(HLM) ratio of variances: intra-individual to sum of inter-individual and inter-family
U	(HLM) scaling factor for correlation $r_z$ , containing measurement error bias

fluctuations in a measurement from one instance to another. Specifically, the SEM formulation assumes one of two possibilities: 1) the absence of measurement errors in the phenotypic data  $y_{i(f(z))}$ , or 2) the inclusion of measurement errors within the nonshared environment component  $E_{i(f(z))}$  (along with its associated variance  $\sigma_E^2$ ) when they are present. Both assumptions lead to the same practical outcome: data with repeated measures are typically preprocessed by aggregating information through a centrality metric (such as the average) before estimating heritability.

Measurement errors are not a significant concern for many measures, such as physical traits that can typically be assessed with high precision. We focus here on heritability estimation in situations where measurement errors are not negligible. For example, in a Stroop task where a large number of trials (e.g., 100) are presented in the experiment for each congruent and incongruent condition. To utilize the SEM formulation directly, the data are typically aggregated, and a centrality measure is used as input.

Data aggregation is a common practice in heritability estimation. Examples can be found in the fields of psychometrics (e.g., Smith et al., 2023; Rea-Sandin et al., 2023; Hung et al., 2023; Gustavson et al., 2023; Vellani et al., 2022; Viktorsson et al., 2022; Yeom et al., 2022; Wang et al., 2020; Routledge et al., 2018; Stins et al., 2005; Schachar et al., 2010; Fan et al., 2001) and neuroimaging (Kastrati et al., 2022; van Drunen et al., 2021; Chen et al., 2019; Harper et al., 2019; van der Meulen et al., 2018; Anokhin et al., 2017; Blokland et al., 2011; Matthews et al., 2007; Polk et al., 2007). However, data aggregation can be a double-edged sword, as valuable information can be lost. The impact of ignoring intra-individual variability has been explored in different contexts. For instance, the issue can be traced back to Spearman (1904) who pointed out the underestimation problem in estimating the correlation between two variables when measurement errors occur. It has also been recently shown that, without proper accountability, test-retest reliability can be substantially underestimated (Rouder and Haaf, 2019; Haines et al., 2020; Chen et al., 2021). Even in heritability estimation, underestimation has been revealed in item response theory due to the adoption of the aggregation process in the form of sum-scores (van den Berg et al., 2007; Schwabe et al., 2019). Recently, the bias problem due to measurement errors has been investigated regarding the reliability of polygenic score and that of a phenotypic trait framed under SEM (Pingault et al., 2022).

Here, we employ a hierarchical linear modeling (HLM) approach to account for intra-individual variability. We propose that *measurement errors should not be considered part of the nonshared environment component*, but rather partitioned appropriately within the model hierarchy. In addition, we employ the HLM approach to reexamine the conventional SEM framework, with the latter being a special case of the former. Specifically, we demonstrate that the common practice of data aggregation fails to adequately address the impact of intra-individual variability. With simulations and real datasets, we address the following questions:

- 1) Does disregarding intra-individual variability result in biased estimation? If so, to what extent?
- 2) To what degree does intra-individual variability contribute to precision in heritability estimation?
- 3) Are typical twin study sample sizes sufficient to achieve a proper precision of heritability estimation?

It is important to note that SEM has been extended to accommodate complex hierarchical data structures (e.g., Mehta and Neale, 2005). Although the SEM framework could potentially be applied to our work, we have opted for HLM due to our preference and familiarity. For clarity, our model comparisons are based on the conventional SEM approach for heritability estimation, not on SEM's extended capabilities.

# 2 Estimating heritability under hierarchical framework

A hierarchical model partitions data variability by mapping the stratified structure to effects at various levels. The modeling approach is well-established in behavior genetics. For example, van den Oord (2001) proposed using HLM to characterize latent genetic and environmental components of variance in extended families. Guo and Wang (2002) derived heritability estimation through direct variance decomposition. McArdle and Prescott (2005) demonstrated the equivalence between conventional SEM and the variance decomposition approach. Hunter (2021) extended this approach to multiple phenotypes, exploring the feasibility of



#### FIGURE 1

Path diagram (reticular action model) for the ACE formulation  $y_i = \alpha$ . All the subscript indices for family (*f*), and zygosity (*z*) are dropped for brevity. Observable effects of *y* are represented as rectangles while latent effects (*A*, *C* and *E* components) are represented by ellipses. A directed path, represented by a single-headed arrow, indicates predictability. An undirected path, represented by a double-headed curved arrow, indicates a covarying relationship. The value on each path shows the correlation coefficient. Note that  $h^2$ ,  $c^2$ , and  $e^2$  can be estimated through effect partitioning into the three latent components (*A*, *C*, and *E*), or through Falconer's formula  $h^2 = 2 (r_{MZ} - r_{DZ})$ .

handling nested data and repeated measures. However, the literature lacks a discussion on how to incorporate intra-individual variability.

### 2.1 HLM: reformulating the SEM approach

We construct the following model by decomposing the data  $y_{i(f(z))}$  into effects through the three indices of individual (i = 1, 2, ..., I), family (f = 1, 2, ..., F), and zygosity (z = MZ, DZ):

individual : 
$$y_i(f(z)) \sim \mathcal{N}(\nu_f(z), \sigma_z^2),$$
  
family :  $\nu_f(z) \sim \mathcal{N}(\alpha, \tau_z^2).$  (7)

Similar to Model 5, the correlation between any two twins,  $i_1$  and  $i_2$ , within a family f can be estimated as

$$r_{z} = \operatorname{corr}\left(y_{i_{1}}(f(z)), y_{i_{2}}(f(z))\right) = \frac{\tau_{z}^{2}}{\tau_{z}^{2} + \sigma_{z}^{2}}, z = MZ, DZ.$$
(8)

The HLM Formulation 7 has been proposed and explored in previous studies in contexts where intra-individual variability is absent (e.g., van den Oord, 2001; Guo and Wang, 2002; McArdle and Prescott, 2005). Here, we aim to extend the framework to accommodate cases where intra-individual variability is present. We highlight three advantages regarding the HLM framework. First, the conventional SEM Formulation 2 assumes the homogeneity of variances (Equation 3), which leads to total variance homogeneity between the two zygosities,

$$\tau_{\rm MZ}^2 + \sigma_{\rm MZ}^2 = \tau_{\rm DZ}^2 + \sigma_{\rm DZ}^2,\tag{9}$$

The equivalence between the two modeling approaches can be established when we equate the total variance of  $\sigma_A^2 + \sigma_C^2 + \sigma_E^2$  in the SEM Formula 2 and that in Equation 9 under the HLM framework (Model 7), as well as  $r_z$  in Models 5, 8, leading to:

$$\begin{split} \sigma_A^2 &= 2 \left( \tau_{MZ}^2 - \tau_{DZ}^2 \right) = 2 \left( \sigma_{DZ}^2 - \sigma_{MZ}^2 \right), \quad \sigma_C^2 = 2 \tau_{DZ}^2 - \tau_{MZ}^2 \\ \sigma_E^2 &= \sigma_{MZ}^2 = \sigma_{DZ}^2 - \left( \tau_{MZ}^2 - \tau_{DZ}^2 \right). \end{split}$$

The HLM (Formulation 7) does not assume variance homogeneity. Instead, a less stringent assumption–proportionality homogeneity across zygosities–is made when using Falconer's Formula 6: the variance proportions accounted for by genetic and shared environmental effects –  $h_z$  and  $c_z$  (z = MZ, DZ) – are the same between the two zygosities (Arbet et al., 2020). Specifically, this proportionality assumption hinges on the derivation (Model 5) for the Falconer's formula. With  $r_z = \rho_z h_z^2 + c_z^2$ , the proportionality assumption of  $h_{MZ}^2 = h_{DZ}^2$  and  $c_{MZ}^2 = c_{DZ}^2$  is sufficient for the validity of the Falconer's formula.

Second, the HLM formulation has the flexibility to accommodate different distribution types. As shown in the Model 7, the individualand family-level variances in the two Gaussian distributions  $\mathcal{N}(\cdot, \sigma_z^2)$  and  $\mathcal{N}(\cdot, \tau_z^2)$  can be expanded to include a wider range of distribution options, such as Student's *t* and exponentially-modified Gaussian distribution. This flexibility in choosing distributions can greatly enhance the quality of the model, especially when working with datasets that exhibit heavy tails, positive-only quantities, truncated or bounded values. Finally, HLM has the ability to explicitly capture intraindividual variability, rather than grouping measurement errors with the nonshared environment component. It allows for appropriately partitioning variability within the data hierarchy.

### 2.2 Consistency between SEM and HLM

We utilized a publicly available dataset of body mass index (BMI) from the R package mets (Holst et al., 2016) to validate the HLM approach. Despite having only one BMI measurement per individual, the presence of intra-individual variability can be considered negligible. In summary, the BMI data comprised I = 11188 twins from F = 6917 families in Finland, including 3665 MZ and 7523 DZ twins aged between 32 and 61 years. Information on each individual's age and sex was also included in the dataset. Both the data and the code for this example are available at https://github.com/afni/apaper\_heritability.

Heritability estimation for the BMI dataset was performed using SEM with the following specifications. Alongside the three latent components A, C, and E in Model 1, we incorporated the following covariates: zygosity and a nonlinear age effect for sex using third-order B-spline bases. The SEM formulation was implemented using the R packages mets and umx (Bates et al., 2019), yielding  $h^2 = 64\%$ ,  $c^2 = 0\%$  and  $e^2 = 34\%$ .

For the HLM approach, we adopted the model (Formulation 7) with log-normal and Gaussian density for the individual- and family-level distributions, respectively based on the tendency of the BMI data skewed to the right (Figure 2A). The following covariates were included: zygosity and nonlinear age effect for each sex using smooth splines with thin plate bases. The model was implemented under the Bayesian framework using the R package brms (Bürkner, 2017). The resulting  $h^2$ ,  $c^2$ , and  $e^2$  were largely consistent with the SEM estimation (Figure 2B), both in terms of point estimate and uncertainty range values.

# 3 HLM: accounting for intra-individual variability

Within the hierarchical framework, we will employ simulations to systematically investigate the influence of intra-individual

variability, trial and participant sample sizes on precision, and determine the necessary sample sizes to achieve a satisfactory level of precision. The insights obtained from these simulations will be further validated by applying the HLM framework to a behavioral dataset.

# 3.1 Measurement errors: part of nonshared environment component?

Measurement errors have traditionally been regarded as part of the nonshared environment component in the heritability model, either implicitly or explicitly (e.g., Maes, 2005; Germine et al., 2015). In other words, it has been considered appropriate to include any measurement errors in the trait measurement  $y_{i(f(z))}$  within the variance component  $\sigma_E^2$  under the SEM (Formulation 2).

We contend that an ideal modeling approach should appropriately allocate measurement errors rather than grouping them together with the nonshared environment. Suppose that the observation  $y_{i(f(z))t}$ (t = 1, 2, ..., T) in the *t*th trial is sampled from a Gaussian distribution with a mean effect  $\theta_{i(f(z))}$  and a standard deviation  $\omega_0$ ,

trial: 
$$y_{i(f(z))t} \sim \mathcal{N}\left(\theta_{i(f(z))}, \omega_{0}^{2}\right).$$
 (10)

Thus, the sample mean  $\hat{\theta}_{i(f(z))} = \sum_{t=1}^{T} y_{i(f(z))t}/T$  carries a crosstrial sampling variance  $\omega^2 = \omega_0^2/T$ , which represents the precision of the estimate. In common practice, when data is aggregated, only the sample mean  $\hat{\theta}_{i(f(z))}$  is utilized in the SEM formulation, while the cross-trial sampling variance  $\omega^2$  is not explicitly accounted for. Consequently,  $\omega^2 = \omega_0^2/T$  remains embedded as part of the nonshared environment component  $\sigma_E^2$ , and the estimation of heritability in Formula 4 becomes dependent on trial sample size and sampling precision. As heritability aims to measure differences among individuals rather than within individuals, it is more conceptually sensible to construct a model where the sample size impacts the precision of the estimated variance, rather than its accuracy.

Measurement errors can be appropriately accounted for as a separate component from the nonshared environment. Suppose we consider the individual-level trait effect  $\theta_{i(f(z))}$  as a latent variable. For the corresponding estimate  $\hat{\theta}_{i(f(z))}$ , we characterize the measurement errors through the cross-trial sampling variance  $\omega^2$ . In other words, we do not solely partition the trait into the three latent components of *A*, *C*, and *E*. Instead, we also treat the true trait effect  $\theta_{i(f(z))}$  as another latent effect, as depicted in the path diagram shown in Figure 3. Additionally, we emphasize that the cross-trial sampling variance  $\omega^2$  is not conceptualized as part of the (latent) nonshared environment component *E*, but rather as a separate entity that is directly observable. Therefore, the original SEM (Formulation 2) is augmented to include two levels,

aggregation : 
$$\hat{\theta}_{i(f(z))} \sim \mathcal{N}\left(\theta_{i(f(z))}, \omega^{2}\right);$$
  
individual :  $\theta_{i(f(z))} = \alpha + A_{i(f(z))} + C_{i(f(z))} + E_{i(f(z))}.$  (11)

The same distribution assumptions in Model 2 apply to the three latent effects of *A*, *C*, and *E* here. Solving this augmented SEM (Formulation 11) directly is challenging. However, in Section 3.4, we will present an approximate approach to heritability estimation.

Measurement errors can also be incorporated into the HLM framework. In particular, we augment the HLM Formulation 7 to

aggregation : 
$$\hat{\theta}_{i(f(z))} \sim \mathcal{N}\left(\theta_{i(f(z))}, \omega^{2}\right);$$
  
individual :  $\theta_{i(f(z))} \sim \mathcal{N}\left(\nu_{f(z)}, \sigma_{z}^{2}\right);$   
family :  $\nu_{f(z)} \sim \mathcal{N}\left(\alpha, \tau_{z}^{2}\right).$  (12)

When the cross-trial sampling variance  $\omega^2$  is available, this augmented model (Formulation 12) can be analyzed under the Bayesian framework or approximately solved as discussed in Section 3.4.

Framing the presence of measurement errors as a distinction between observed and latent effects helps in understanding the associated impact. As illustrated in the path diagram (Figure 3), one can estimate heritability directly based on the correlations  $\hat{r}_z$  using the sample means  $\hat{\theta}_{i(f(z))}$  through Falconer's Formula 6. Long ago, Spearman (1904) highlighted a bias issue: the correlation between two quantities becomes attenuated when measurement errors are not accounted for. Similarly, when measurement errors are disregarded,  $\hat{r}_z$  would be underestimated compared to their counterparts  $r_z$  based on the latent effects  $\theta_{i(f(z))}$  (Figure 3). Next, we construct an HLM formulation at the observation level to fully illustrate the issues associated with data aggregation.

## 3.2 HLM with trial-level data under one task condition

We begin by extending the HLM (Equation 7) to a case where data are collected at the observation level with repeated measures through trials under a single task condition. The data  $y_{i(f(z))t}$  are represented using four indices: family (f = 1, 2, ..., F), zygosity (z = MZ, DZ), individual (i = 1, 2, ..., I), and trial (t = 1, 2, ..., T). Instead of utilizing aggregated information as in Equations 11, 12, we formulate the following HLM based on Formulation 10:

trial: 
$$y_i(f(z))_t \sim \mathcal{N}\left(\theta_i(f(z)), \omega_0^2\right);$$
  
individual:  $\theta_i(f(z)) \sim \mathcal{N}\left(v_f(z), \sigma_z^2\right);$   
family:  $v_f(z) \sim \mathcal{N}\left(\alpha, \tau_z^2\right).$  (13)

The corresponding path diagram is shown in Figure 4. The Model 13 can be extended to include various distributions. For instance, the trial-level effects can be characterized through Bernoulli, gamma, and Poisson distributions, which allow for modeling different types of data (count, binary, skewed).

Heritability can be estimated using the HLM (Formulation 13) for trial-level data as follows. Similar to the situation with HLM for the conventional SEM in (Formulation 7), we compute the correlations between two twins within a family using the inter-individual and interfamily variances,  $\sigma_z^2$  and  $\tau_z^2$ , through the Formula 8. Then, the three variance proportions  $h^2$ ,  $c^2$ , and  $e^2$  are obtained using Falconer's Formula 6.

# 3.3 Consequences of data aggregation under the SEM formulation

Now we examine the common practice of data aggregation in light of the HLM (Formulation 13). To capture the crucial role of



(A) Data distribution. BMI exhibits a right-skewed distribution, with slightly greater dispersion among D2 twins compared to M2 twins. (B) Heritability estimation of BMI data. The proportions of BMI variability attributed to each of the three components are depicted. The distributions for  $h^2$ ,  $c^2$ , and  $e^2$  are estimated using the HLM Formulation 7 and presented in blue. The shaded area under each distribution represents the 95% highest density interval, while the vertical dashed line indicates the mode (peak). For comparison, the point estimates (dot) and their corresponding 95% uncertainty intervals (horizontal line) for the SEM (Formulation 2) are displayed in red.

intra-individual variability  $\omega_0^2$  across different phenotypic traits, we define a dimensionless measure of the variability ratio for each zygosity:

$$R_{\nu,z} = \sqrt{\frac{\omega_0^2}{\tau_z^2 + \sigma_z^2}}, z = MZ, DZ.$$
 (14)

The variability ratio  $R_{\nu,z}$  captures the fundamental aspect of heritability: the proportion of inter-family and inter-individual variance relative to intra-individual variance. Under the assumption of homogeneity (Equation 9),  $R_{\nu,MZ} = R_{\nu,DZ}$ . For simplicity, we drop the subscript Z and denote their average as  $R_{\nu}$ .

When the trial-level data  $y_{i(f(z))t}$  are aggregated across trials with their average  $\bar{y}_{i(f(z))} = \frac{1}{T} \sum_{t=1}^{T} y_{i(f(z))t}$ , the model (Formulation 7) becomes

$$\begin{aligned} \text{individual} : \bar{y}_{i(f(z))} \sim \mathcal{N}(\nu_{f(z)}, \tilde{\sigma}_{z}^{2}), \ \tilde{\sigma}_{z}^{2} = \sigma_{z}^{2} + \omega_{0}^{2} / T; \\ \text{family} : \nu_{f(z)} \sim \mathcal{N}(\alpha, \tau_{z}^{2}). \end{aligned}$$
(15)

In comparing the model (Formulation 15) with aggregated data to its counterpart (Formulation 7) for data without measurement errors, we note that ignoring intra-individual variability leads to its combination with the inter-individual variance  $\sigma_z^2$  into  $\tilde{\sigma}_z^2$ . As a result, the correlations  $r_{\rm MZ}$  and  $r_{\rm DZ}$  in Formula 8 are updated to

$$\tilde{r}_{z} = \frac{\tau_{z}^{2}}{\tau_{z}^{2} + \tilde{\sigma}_{z}^{2}} = \frac{\tau_{z}^{2}}{\tau_{z}^{2} + \sigma_{z}^{2} + \omega_{0}^{2}/T} = r_{z}U,$$
(16)

where the introduced bias into  $r_{MZ}$  and  $r_{DZ}$  is quantified by the dimensionless quantity  $U = \frac{1}{1+R_v^2/T}$ . It is noteworthy that when  $\omega_0$  is nonzero, U < 1, signifying that  $\tilde{r}_z$  is consistently downward biased. This bias arises due to the presence of intra-individual variability, and its attenuation rate follows a sigmoid function of  $R_v$ . In the limit where  $\omega_0^2/[T(\tau_z^2 + \sigma_z^2)] \rightarrow 0$ , which can occur with decreasing standard error  $\omega_0$  or increasing trial size T,  $\tilde{r}_z \rightarrow r_z$ .

The parameter U quantifies the degree of bias in heritability estimation under SEM when data aggregation is applied. According to Falconer's Formula 6, both  $h^2$  and  $c^2$  would be underestimated by a factor of U, while  $e^2$  would be



overestimated by the same factor. For example, with T = 100 trials, a small intra-individual variability such as  $R_v = 1$  has negligible impact on heritability estimation ( $U \approx 0.99$ ), whereas a large intra-individual variability with  $R_v = 10$  substantially underestimates  $h^2$  and  $c^2$  by 50%. Conversely, if  $R_v = 3$ , biases cannot be disregarded even with T = 20 trials unless T approaches or exceeds 100.

One direct way to view the distinction between SEM and HLM is to compare their respective path diagrams (Figures 3, 4). HLM preserves the hierarchical structure and cross-trial variability, ensuring this information propagates across other hierarchical levels (Model 13). In contrast, SEM obscures this variability through data aggregation, compromising the hierarchical integrity. This loss of data structure fidelity in SEM leads to biased underestimation of heritability, as captured through the parameter U in the expression (Formula 16).

# 3.4 Ameliorating the biases in the SEM formulation

The biases induced in the SEM formulation can be theoretically corrected by introducing an adjustment term in the denominator of (Formula 16) to counteract the contaminating effect of  $\omega_0^2/T$  under the model (Formulation 15),

$$r_z = \frac{\tau_z^2}{\tau_z^2 + \tilde{\sigma}_z^2 - \omega_0^2 / T}.$$
(17)

Similarly, decontamination can be achieved for the Formulas 4, 8. However, these adjustments rely on the availability of the intraindividual variance  $\omega_0^2$ , which is not directly accessible once the data are aggregated. Nevertheless, the biases can be practically mitigated. For example, we can use the cross-trial variance estimates  $\hat{\omega}_0^2 = \frac{1}{I} \sum_{i=1}^{I} s_{i(f(z))}^2$  in Formula 17, leading to the following approximate adjustment,

$$\hat{r}_{z} = \frac{\tau_{z}^{2}}{\tau_{z}^{2} + \tilde{\sigma}_{z}^{2} - \hat{\omega}_{0}^{2} / T}.$$
(18)

Similarly, as  $\sigma_E$  inherently contains the additive contribution of measurement error, we can directly adjust the biases in the SEM estimates to,

$$\hat{h}^{2} = \frac{\sigma_{A}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2} - \hat{\omega}_{0}^{2}/T}, \quad \hat{c}^{2} = \frac{\sigma_{C}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2} - \hat{\omega}_{0}^{2}/T}, \quad \hat{e}^{2}$$
$$= \frac{\sigma_{E}^{2} - \frac{1}{T}\hat{\omega}_{0}^{2}}{\sigma_{A}^{2} + \sigma_{C}^{2} + \sigma_{E}^{2} - \hat{\omega}_{0}^{2}/T} = 1 - \hat{h}^{2} - \hat{c}^{2}. \tag{19}$$

These approximate adjustments (Formulas 18, 19) offer a solution to the augmented SEM Formulation 11 and its hierarchical counterpart (Model 12). We will further explore and validate their effectiveness as approximate adjustments later with an experimental dataset.

An intriguing aspect of biased estimation for heritability is its analogy to the phenomenon of correlation attenuation in the presence of measurement errors. Spearman (1904) recognized the problem of bias caused by measurement errors and proposed an adjustment method to disattenuate the correlation between two variables. In Formula 16, the term U serves a similar purpose to the reliability coefficient or separation index in classical test theory. Consequently, it is interesting to note that the decontamination Formula 18 and its approximation (Formula 19) employ a similar adjustment strategy as suggested by Spearman (1904).

## 3.5 HLM with trial-level data under two task conditions

We now extend the HLM (Formulation 13) to accommodate two task conditions. In fields such as psychometrics and neuroimaging, the focus often lies in comparing and analyzing the contrast between two conditions. We expand the previous HLM (Formulation 13) to one with hierarchical levels using five indices: family (f = 1, 2, ..., F), zygosity (z = MZ, DZ), individual (i = 1, 2, ..., I), condition ( $c = c_1, c_2$ ), and trial (t = 1, 2, ..., T):



**FIGURE 4** Path diagram (reticular action model) for the HLM Formulation 13. Subscript indices for family *f*, and zygosity *z* are omitted from the nodes for brevity. Unlike the path diagram in Figure 3, the trait measures *y* are observable. A directed path, indicated by a singleheaded arrow, represents predictability, while an undirected path, shown by a double-headed curved arrow, represents a covarying relationship. The value on each path indicates the correlation coefficient.

$$\text{trial}: y_{ci}(f(z))_{t} \sim \mathcal{N}\left(\theta_{ci}(f(z)), \omega_{0}^{2}\right);$$
  

$$\text{ndividual}: \begin{bmatrix} \theta_{c_{1}i}(f(z)) \\ \theta_{c_{2}i}(f(z)) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \nu_{c_{1}f(z)} \\ \nu_{c_{2}f(z)} \end{bmatrix}, \begin{bmatrix} \sigma_{c_{1},g}^{2} & \ast \\ \ast & \sigma_{c_{2},g}^{2} \end{bmatrix}\right);$$
  

$$\text{family}: \begin{bmatrix} \nu_{c_{1}f(z)} \\ \nu_{c_{2}f(z)} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \alpha_{c_{1}} \\ \alpha_{c_{2}} \end{bmatrix}, \begin{bmatrix} \tau_{c_{1},g}^{2} & \ast \\ \ast & \tau_{c_{2},g}^{2} \end{bmatrix}\right).$$
(20)

iı

The differences from Section 2.2 are twofold: (a) the presence of two intercepts,  $\alpha_c$ , one for each condition, and (b) the individual- and family-level distributions being bivariate instead of univariate. While the covariances for the individual- and family-level distributions are not of interest in the current context, we acknowledge their presence by denoting them with an asterisk \* in the respective variance-covariance matrix.

Heritability estimation is straightforward for each condition under the HLM (Formulation 20). First, we calculate the correlation between two twins within a family for condition  $c_k$  (k = 1, 2) using

$$r_{c_k,MZ} = \frac{\tau_{c_k,MZ}^2}{\tau_{c_k,MZ}^2 + \sigma_{c_k,MZ}^2},$$

$$r_{c_k,DZ} = \frac{\tau_{c_k,DZ}^2}{\tau_{c_k,DZ}^2 + \sigma_{c_k,DZ}^2}.$$
(21)

Then, we estimate heritability for each condition by plugging  $r_{c_k,MZ}$  and  $r_{c_k,DZ}$  into Falconer's Formula 6.

Two approaches are available for estimating heritability and the variability ratio for the contrast between two conditions. One approach involves reparameterizing the model (Formulation 20) through an indicator variable using effect coding for the two conditions,

$$x_{c_k} = \begin{cases} 0.5, & \text{if } k = 1; \\ -0.5, & \text{if } k = 2. \end{cases}$$
(22)

Alternatively, within the Bayesian framework, one can directly obtain the distribution of the contrast by formulating it based on each condition's posterior draws from the HLM (Formulation 20). Biased estimation due to data aggregation and its adjustment in the previous subsection also apply to the case with two conditions. For each condition and their contrast, the variability ratio  $R_v$  can be similarly defined. The only modification for the contrast is to replace  $\sigma_0^2/T$  with  $2\sigma_0^2/T$  in Formula 14. Similarly, the discussions regarding the biases in heritability estimation for each condition in the preceding subsection can be directly applied here. However, when considering the contrast, the bias requires replacing  $\sigma_0^2/T$ in (Formulas 15–19) with  $2\sigma_0^2/T$ .

# 3.6 Assessing model performance through simulations

We have shown that, by closely reflecting the actual data structure, an HLM framework appropriately accounts for different sources of data variability. However, the precision of heritability estimation remains unclear, as there is no analytical quantification available. Simulations were conducted to assess the precision of heritability estimation from the following aspects regarding the impact of intra-individual variability that the analytical approach cannot easily reveal: 1) the precision of heritability estimation, and 2) the requirement of family and trial sample sizes in twin studies.

Detailed information about the simulations and results can be found in Supplementary Appendix A. As expected, HLM showed no bias in heritability estimation, whereas biases under the SEM framework become more pronounced as the variability ratio  $R_{\nu}$ increases and/or the trial sample size *T* decreases. Additionally, biases under the SEM formulation with aggregated data can be effectively adjusted using empirical or theoretical standard errors through Formula 17 or Formula 18. More importantly, simulations indicate that intra-individual variability impacts estimation precision. Specifically, larger  $R_{\nu}$  leads to poorer precision. Lastly, family sample size has a greater impact than trial sample size on estimation precision. For instance, when  $R_{\nu} \leq 1$ , an appropriate level of uncertainty can be achieved with 50 trials and 1,000 families. When  $R_{\nu} \gg 1$ , several thousand families may be necessary.

# 4 Applying HLM to an experimental dataset

We apply the HLM approach to an experimental dataset to address two primary questions. First, do the insights gained from numerical simulations in the previous section align with the findings when real data is analyzed? Second, what is the range of the relative magnitude of intra-individual variability, as indicated by the ratio  $R_{\nu}$ , in commonly encountered empirical datasets?

We utilize a behavioral dataset obtained from an experiment conducted as part of the ABCD study. The experiment investigated selective attention during adolescence using an emotional Stroop task (Smolker et al., 2022), with reaction time (RT) as a phenotypic trait. The data was collected during the 1-year follow-up visit and is publicly accessible through the 2020 ABCD Annual Curated Data Release 4.0 (https://nda.nih.gov/study.html?id=1299). The analysis scripts can be found online at: https://github.com/afni/apaper\_heritability.

### 4.1 Data description

A subset of the original dataset, specifically containing twins, was utilized for the analysis. Refer to Table 2 for detailed information on participant counts and demographic data. The subset consisted of 1,102 twins (including some triplets) from 555 families and was selected from a larger dataset of 11,876 participants (Iacono et al., 2018). Among the included twins, there were 461 MZ and 641 DZ participants.

We focus on the RT data for estimating heritability. The RT was measured for two levels of congruency in a Stroop task: congruent and incongruent. Each participant was instructed to respond to a total of 48 trials, consisting of 24 congruent and 24 incongruent trials. The response window for each trial was set to 2,000 ms. Among 1,102 twins, a total of 49,524 trials were included in the analysis after excluding incorrect responses. This resulted in an overall correct response rate of approximately 93.6%. The distribution of RTs is heavily right-skewed (Figure 5), with a mode of 962 ms and a 95% highest density interval ranging from 634 to 1,827 ms.

### 4.2 Model comparisons with real data

The RT data was analyzed using two different approaches: HLM and SEM. For HLM, the trial-level data was fitted using the Formulation 20. Site, sex, race, age, and zygosity were included as covariates. To account for the right-skewness (Figure 5), a lognormal distribution was used for the trial-level effects. The Bayesian framework was employed to implement the HLM approach, utilizing the brms package in R (Bürkner, 2017). Heritabilty was estimated for each of the three RT effects: the congruent condition, incongruent condition, and the Stroop effect (the contrast between incongruent and congruent conditions). The SEM was applied with aggregated data. Specifically, RT was aggregated across trials for each condition at the individual level. Similar to the HLM approach, covariates including site, sex, race, age, and zygosity were included. The SEM formulation was implemented using the R packages of mets and umx. The computational time for SEM with aggregated data was negligible. In contrast, HLM-based estimation, using Markov chain Monte Carlo simulations, required 3.5 hours with four chains and 12 threads per chain. These computations were performed on an Intel Server S2600WFT equipped with 96 CPUs running at 2,933 MHz.

The estimation results are presented in Figure 6. The performance of HLM relative to SEM, as summarized below, aligns with our simulations in Section 3.6. Overall, both SEM and HLM exhibited a significant amount of uncertainty in estimating heritability for both conditions. SEM displayed noticeable underestimation of  $h^2$  and  $c^2$  (first two columns, Figure 6). However, neither modeling approach provided satisfactory estimation for the contrast (third column, Figure 6). This estimation challenge arises from the combination of three factors encapsulated by a large intra-individual variability ratio  $R_v \approx 10$ : 1) a much smaller effect size, 2) an extremely limited trial sample size, and 3) a relatively small twin sample size.

Below are a few detailed elaborations:

- 1) Impact of relative intra-individual variability on estimation precision. First, under each of the individual congruent and incongruent conditions, the observed intra-individual variability was not large  $(R_v \approx 2)$ , resulting in moderate uncertainties for HLM estimates of  $h^2$ ,  $c^2$ , and  $e^2$ . However, the SEM approach showed difficulty in accurately assessing uncertainty near the parameter boundaries (e.g., 0 or 1 for  $h^2$ ,  $c^2$ , and  $e^2$ ). For instance, the SEM's small uncertainty interval (0, 0.01) for  $c^2$  likely stemmed from numerical singularity issues in the traditional statistical framework. In contrast, regularization in hierarchical modeling (Chung et al., 2013) yielded a more reasonable uncertainty interval (0, 0.28) for  $c^2$ under the Bayesian framework. Second, the intra-individual variability for the contrast between the two conditions was large ( $R_v \approx 10$ ). Consequently, the uncertainties for  $h^2$ ,  $c^2$  and  $e^2$  were very large, with the estimated density of  $c^2$  resembling a uniform distribution. The SEM approach did not provide any meaningful estimates either, and this challenge was further demonstrated by its inability to provide an appropriate uncertainty interval, yielding only a single point estimate constrained at the parameter boundaries, likely due to convergence difficulties.
- 2) Impact of relative intra-individual variability on estimation bias. Under both conditions, the intra-individual variability is moderate ( $R_v \approx 2$ ), resulting in small underestimations of  $h^2$ and  $c^2$  by SEM. The overestimation of  $e^2$  was also small. However, the large intra-individual variability for the contrast ( $R_v \approx 10$ ) led to more noticeable underestimations of  $h^2$  and  $c^2$  by SEM.
- 3) Impact of relative intra-individual variability on sample sizes. The larger  $R_{\nu}$  for the contrast ( $R_{\nu} \approx 10$ ) is consequential. Simulation results in Section 3.6 indicate that larger sample sizes, especially in terms of family count, would be required to reduce the large uncertainty. We note that the observed range of  $R_{\nu}$  values aligns with psychometric data from individual studies in test-retest reliability estimation (Rouder and Haaf, 2019; Chen et al., 2021; Baker et al., 2021).
- 4) Bias adjustment for SEM estimates. The biases under the SEM framework, due to data aggregation, adjusted using Formula 18, were reduced. The adjusted estimates for  $\hat{h}^2$  under the congruent, incongruent, and contrast conditions were 0.38, 0.40, and 0.0, respectively (green triangles, Figure 6). These adjustments effectively reduced bias, although they remained slightly biased compared to HLM, which could be attributed to deviations from the Gaussian assumption under SEM.

We also explored the HLM approach for the ABCD-Stroop data using a conventional linear mixed-effects modeling framework instead of a Bayesian approach. The lmer function from the lme4 package in R (Bates et al., 2015) was utilized to fit the models (Formulations 13, 20) with RT log-transformed. Although the point estimates (not shown here) for  $h^2$ ,  $c^2$ ,  $e^2$ , and  $R_v$  under the congruent and incongruent task conditions were largely consistent with the values obtained under the Bayesian framework, the numerical solver in lme4 failed to converge for the contrast between the incongruent and congruent conditions due to the relatively small inter-individual variances ( $R_v \approx 10$ ).

## 5 Discussion

Heritability estimation based on data with a non-negligible intra-individual variability necessitates a model that accurately represents the underlying hierarchical structure. When assessing a phenotypic trait with repeated measures, we propose a hierarchical model that encompasses all relevant levels, allowing for the incorporation, propagation, and separation of intra-individual measurement error from parameter estimation at higher levels. Through numerical simulations and a real behavioral dataset, we have demonstrated a few advantages of HLM. These advantages include: 1) avoidance of estimation bias, 2) the ability to account for the significant influence of intra-individual variability on heritability estimation, and 3) enhanced interpretability and explanatory power of results, such as identifying the challenges associated with reducing estimation uncertainty due to sample size limitations.

# 5.1 The importance of modeling data generating mechanism

Data reduction through aggregation is a commonly used in statistical applications, particularly in studying individual differences. Even though intraindividual variability has long been recognized (Fiske and Rice, 1955), many classical frameworks have been applied in contexts where measurement errors are minimal or nonexistent. For example, intraclass correlation (Fisher, 1954) for test-retest reliability in individual differences is typically assessed without considering intra-individual variability. A similar situation can be observed in heritability estimation: to date, common modeling approaches do not explicitly account for the level of measurement units (e.g., individual trials), and instead they simplify the data through preliminary aggregation steps such as averaging.

Can data aggregation be justified by attributing intraindividual variability to the nonshared environment (the *E* component in the ACE model)? The underlying rationale for data aggregation is that intraindividual variability arises either from true biological fluctuations (ontological variability) or measurement limitations (epistemological noise). However, the specific sources are often too complex to be fully accounted for in typical studies. A more pragmatic and effective approach is to adopt a causal inference perspective that focuses on the underlying data-generating process. It is well established that when within-individual variability follows systematic patterns, treating it solely as residual errors can introduce bias and misinterpretation. Instead, explicitly modeling this variability is crucial to ensuring accurate and meaningful estimates.

In heritability estimation, path diagrams are commonly used to depict causal relationships among variables. Within this framework, a latent trait or condition is conceptualized as a higher-level theoretical construct that causally influences each individual measurement or trial (see Figure 4). In other words, individual measurements are specific realizations determined by the underlying latent construct. Crucially, heritability is defined at the level of this latent construct, not at the level of single measurements.

A hierarchical modeling framework more accurately maps the causal structure outlined in the path diagram. In contrast, conventional SEM, which treats intraindividual variability as

TABLE 2 Demographic information of twins in a Stroop experiment from the ABCD Study.

Twin	I = 1102 twins; $F = 555$ families; 3 families with DZ triplets; 11 families with available data from only one twin
Zygosity	461 MZ twins, 641 DZ twins
Sex	550 males, 552 females
Race	720 white, 152 black, 112 Hispanic, 3 Asian, 114 others
Age	108-132 months; mean: 121 months, standard deviation: 6 months



FIGURE 5

Reaction time (RT) distributions. The shaded area under each density represents the 95% highest density interval, while the vertical dashed line indicates the mode. The RT distribution is 1) right-skewed, 2) slightly right-shifted under the incongruent condition compared to the congruent condition, and 3) slightly more dispersed for DZ twins than MZ twins.

residual errors, does not fully align with the causal structure and can lead to underestimated heritability, as demonstrated in Section 3.3. This underscores the necessity of explicitly modeling hierarchical data structures to account for both between- and within-individual variability.

Our empirical evidence here from both a real dataset and simulations demonstrates that adopting an HLM framework-which respects the full hierarchical structure of the data-yields more accurate heritability estimates than the conventional SEM for psychometric traits. When intraindividual variability is minimal-as is often the case with many physical traits-the conventional SEM can be seen as a special asymptotic case of a more general HLM, and aggregation is justified because the impact of such variability is negligible ( $R_v \ll 1$ ). However, for traits such as psychometric measures where intraindividual variability is substantial  $(R_v \gg 1)$ , simply aggregating data (i.e., incorporating this variability into residual errors) will likely lead to biased heritability estimates. This underestimation issue extends beyond heritability estimation and has also been observed in test-retest reliability estimation (Rouder and Haaf, 2019; Haines et al., 2020; Chen et al., 2021) and in neuroimaging experimental designs, where the role of trial samples is often overlooked (Chen et al., 2022).

To recapitulate, the HLM approach acknowledges that heritability is defined at the latent trait level and avoids biases associated with oversimplified data aggregation. While data aggregation may be acceptable for traits with minimal intraindividual variability, a hierarchical modeling approach that directly incorporates the causal structure of the data is essential for accurately estimating heritability when variability is pronounced. To improve the accuracy of heritability estimation, we recommend adopting HLM in the presence of intra-individual variability. It is important to note that the HLM framework is not mutually exclusive with SEM, on which SEM and other methods such as common pathway model are based. As the path diagrams in Figures 1, 4 illustrate, both frameworks are conceptually consistent, as discussed in Section 2.2. Nevertheless, we emphasize that the broader framework of HLM combined with the Bayesian approach offers several advantages:

- 1) It supports a wider range of numerically implemented distributions (e.g., Student's *t*, inverse Gaussian).
- 2) It integrates uncertainty assessment into a single process.
- 3) It robustly handles variance-covariance structures.

While the last point is important for theoretical and interpretational reasons, it also has useful practical benefits. In the commonly-used R software packages, there are severe challenges faced when using methods implemented in the nlme and lme4 packages, which can struggle with numerical singularities when correlations (or variances) approach boundary values such as -1, 1, or 0, as encountered in this Stroop dataset. The proposed framework avoids these difficulties.

# 5.2 Biases, uncertainty and challenge of heritability estimation

There are two aspects of accuracy compromise that need to be considered in heritability estimation. The first aspect pertains to estimation biases. As demonstrated in this study, failure to fully incorporate the data structure can lead to biased estimates of heritability. The second aspect concerns the uncertainty in heritability estimation. In addition to providing a point estimate for the effect of interest, it is equally important to quantify its uncertainty, characterized through measures such as standard error, an uncertainty interval (e.g., 95%), or even a full distribution (as depicted in Figure 6). However, uncertainty is often not well emphasized in common practice. In some cases, only the central tendency (e.g., mean) of heritability estimation is reported. However, to truly comprehend the generalizability of results, understanding uncertainty is crucial. One of the benefits of the HLM framework is its ability to directly generate posterior distributions that illustrate estimate precision.

In the presence of intra-individual variability, one might be tempted to adopt the bias adjustment approach using the conventional SEM. Our findings demonstrate that the biases resulting from data aggregation can be mitigated to some extent if variability can be determined separately (e.g., through repeated measures), as indicated by Formula 18 or Formula 19. However, in practical applications, these adjustments are suboptimal due to distributional deviations, as demonstrated in our example using the Stroop dataset. Furthermore, an effective adjustment for biases in uncertainty assessment is currently lacking. Hence, a comprehensive HLM framework remains the preferred choice.

Sample sizes remain a challenge in twin studies. The dataset we used for demonstration, exemplifying a cognitive inhibition study, suggests that reasonable levels of uncertainty can be achieved with



### FIGURE 6

Estimated heritability for RT data. (A) The three columns represent the two conditions (congruent and incongruent), as well as their contrast, with the corresponding variability ratio  $R_v$  indicated in the column labels. The three rows correspond to  $h^2$ ,  $c^2$ , and  $e^2$ . In each panel, the HLM result is represented by a solid blue density curve, derived from random draws from posterior chains. The mode is marked by a vertical blue dashed line, and the shaded blue region represents the 95% uncertainty interval. The SEM counterparts are also displayed in each plot, with the point estimate depicted as a red dot and its 95% uncertainty interval represented by a horizontal red line. Notably, the SEM point estimates tend to be smaller than their HLM counterparts for  $h^2$  and  $c^2$  (while larger for  $e^2$ ). Adjustments for the SEM estimates, using the Formula 19, are denoted as SEM1 and shown as green triangles. (B) Comparisons among the three models are presented with their point estimates and 95% uncertainty intervals. For the HLM, estimates are derived from the modes and highest density intervals of the posterior distributions in (A).

sample sizes of less than 1,000 families and less than 100 trials per individual condition (congruent and incongruent). The estimated heritability of approximately 40% (first two columns, Figure 6) aligns with the general range observed in typical phenotypic traits in the literature (Polderman et al., 2015). However, the contrast between conditions is often the focal point of interest. Even with HLM estimation, the uncertainty of heritability for this contrast remains unresolved (third column, Figure 6), creating imprecision regarding its magnitude. In other words, despite attempts by the Consortium (Iacono et al., 2018; Smolker et al., 2022) to address the sample size issue, the dataset from the ABCD Study (consisting of 461 MZ and 641 DZ twin pairs, with less than 48 trials per condition) does not provide sufficient certainty for estimating the heritability of the Stroop effect. Achieving a reasonable level of precision may require impractical sample sizes (e.g., hundreds of trials and thousands of individuals).

The relative magnitude of intra-individual variability, as quantified by the ratio  $R_{\nu}$ , serves as an informative indicator in heritability modeling. As a dimensionless parameter, it influences not only the accuracy and uncertainty of heritability estimation but also those of test-retest reliability (Rouder et al., 2019; Chen et al., 2021). Historical power analyses in twin studies have suggested a minimum sample size of 600 twin pairs (Martin et al., 1978; Sham

et al., 2020). However, our simulations demonstrate that, in the presence of measurement errors, a large  $R_{\nu}$  poses a significant challenge for future studies in the field of individual differences, particularly when examining effect contrasts and higher-order interactions. Additionally, this ratio highlights the relative importance of trial sample size compared to participant sample size across various experimental modalities, such as functional magnetic resonance imaging, magnetoencephalography, electroencephalography, and psychometrics. In all these cases, the  $R_{\nu}$  ratio exceeds 1, and sometimes even surpasses 10 (Baker et al., 2021; Chen et al., 2021; Chen et al., 2022). Due to this substantial ratio, the sample size of trials can be nearly as crucial as the number of participants in terms of experimental efficiency in neuroimaging and psychometrics.

### 5.3 Heritability estimation in neuroimaging

To date, there has been an increasing number of twin studies utilizing task-based functional magnetic resonance imaging (fMRI). In these studies, it has been a common practice to aggregate data across trials during fMRI data analysis, resulting in the neglect of intra-individual variability, which is neither accounted for nor reported. For instance, Polk et al. (2007)

reported a small heritability estimate of blood oxygenation leveldependent (BOLD) response  $(h^2 \sim 0.2)$  for face and house processing (the specific contrasts were not analyzed), but negligible heritability for pseudowords and chairs in the ventral visual cortex, based on an fMRI experiment involving 13 MZ and 11 DZ twins, with 90 trials per condition. Similarly, Matthews et al. (2007) revealed a moderate heritability  $(h^2 = 0.37; 90\%$  interval: (0, 0.74)) for the interference effect in the dorsal anterior cingulate cortex during a multi-source interference task with congruent and incongruent conditions, involving 20 MZ and 20 DZ twins, with 144 trials per condition. The heritability estimates for other regions were negligible. In contrast, the heritability of reaction time was moderate for the congruent condition ( $h^2 = 0.45$ ; 90% interval: (0, 0.76)), but negligible for the incongruent condition and the interference effect. Additionally, Blokland et al. (2011) found moderate to high heritability estimates ( $h^2 = 0.40$  to 0.65) in more than ten regions during an n-back working memory experiment involving 150 MZ and 132 DZ twins, with 128 trials per condition.

Is intra-individual variability a concern for heritability estimation in neuroimaging? The aforementioned task-based fMRI experiments have primarily relied on a large number of trials to obtain reliable estimates of condition-level effects. This is similar to the Stroop dataset we investigated here, with the distinction that the focus is on BOLD response rather than reaction time. However, the family sample size has often been relatively small, leading to larger uncertainty ranges in the estimates. This issue is particularly pronounced because the relative intra-individual variability across the brain, as reported in the literature, tends to be substantial, with  $R_v \gg 1$  (Chen et al., 2021; Baker et al., 2021; Chen et al., 2022). Thus, heritability estimation in neuroimaging is at least as challenging as typical traits such as psychometric data.

# 5.4 Limitations of heritability estimation through HLM

The HLM approach comes with additional costs. Firstly, introducing an extra level in the data hierarchy significantly increases the complexity of the model structure. Secondly, and perhaps the greater challenge, this increased complexity brings along numerical burdens. Traditional tools like linear mixed-effects estimation are not well-suited for solving hierarchical models of heritability. Instead, resorting to a Bayesian approach may be necessary to handle the numerical challenges (e.g., singularity).

There is always room for improvement in modeling. For example, the full details of the underlying mechanism and framework of cognitive inhibition involved in the reaction time of the Stroop effect are not fully known to researchers. Therefore, no model can fully replicate their structure. However, HLM attempts to model as much as is known and observed in a study. Model fitting can be improved by incorporating auxiliary information, such as accommodating abnormalities like skewness, outliers, and truncation through more inclusive and adaptive distributions (e.g., log-normal, ex-Gaussian). Additionally, one could reconsider the chosen partitioning into three components of  $h^2$ ,  $c^2$ , and  $e^2$  in twin studies and other assumptions (Robette et al., 2022): the additivity of genetic effects, the absence of assortative mating, the nonexistence of genetic dominance or epistasis, the generalizability from twins to the rest of the population, equal environment impact between MZ and DZ twins, and the absence of gene-environment correlation or interaction.

Further integrating HLM with the conventional SEM framework presents a promising avenue for future research. SEM, with its long-established history, offers distinct advantages, including intuitive interpretation, specialized applications, and computational efficiency. While beyond the scope of this study, leveraging the strengths of both SEM and HLM (e.g., Mehta and Neale, 2005) holds significant potential. A unified approach could enable greater flexibility in modeling distributions, account for intraindividual variability, and reduce estimation biases. In addition, this study focuses on within-individual categorical variables, such as task conditions. Future research could extend this framework to within-individual quantitative variables, particularly in the context of longitudinal data (e.g., Eaves et al., 1986).

The interpretation of heritability is subtle and sometimes controversial. Our focus here is solely on the technical aspects of heritability estimation. Nevertheless, we emphasize caution in its interpretation. As a statistical metric, heritability captures variation and/or correlation rather than causation. Therefore, one must not confuse the extent of phenotype variability with the contribution of genetic factors. The concept of heritability effectively pertains to the population level and cannot be realistically applied to a particular individual. On the other hand, the information provided by heritability lies in its potential predictability. It can probabilistically predict, but not causally determine, the extent of phenotypic variability. A high heritability for a phenotypic trait may warrant further investigation into the underlying complex genetic mechanisms or the etiology of genetic risk factors, such as biomarkers. This perspective highlights the need to complement heritability research of variance partitioning with mechanism elucidation (Downes and Turkheimer, 2022).

## 6 Conclusion

We propose an HLM approach to improve heritability estimation in twin studies when the phenotypic trait is measured with multiple samples. The methodology aims to separate measurement errors from the variations of interest and addresses issues such as information loss due to data reduction, distribution violations, and uncertainty characterization in current modeling approaches. We demonstrated that the conventional SEM is likely to underestimate heritability when intraindividual variability is moderate to high (which is common in many real-world scenarios). We supported this finding with analytical derivations, simulations and an experimental dataset from the ABCD study, validating the performance of the HLM approach. Our simulation results suggest that traits with small effect sizes may require much larger sample sizes than currently practiced.

## Data availability statement

The data presented in the study are deposited in the repository https://nda.nih.gov/study.html?id=1299, accession number 1299.

### **Ethics statement**

The studies involving humans were approved by Institutional Review Board of the University of California, San Diego (UCSD). The studies were conducted in accordance with the local legislation and institutional requirements. Written informed consent for participation in this study was provided by the participants' legal guardians/next of kin.

### Author contributions

GC: Conceptualization, Data curation, Formal Analysis, Investigation, Methodology, Project administration, Resources, Software, Supervision, Validation, Visualization, Writing–original draft, Writing–review and editing. DM: Data curation, Writing–review and editing. PT: Funding acquisition, Supervision, Writing–review and editing.

### Funding

The author(s) declare that financial support was received for the research, authorship, and/or publication of this article. GC and PT were supported by the NIMH Intramural Research Program (ZICMH002888) of the NIH/HHS, United States. DM was also supported by the NIMH Intramural Research Program (ZICMH002960) of the NIH/HHS, United States. Data used in

### References

Anokhin, A. P., Golosheykin, S., Grant, J. D., and Heath, A. C. (2017). Heritability of brain activity related to response inhibition: a longitudinal genetic study in adolescent twins. *Int. J. Psychophysiol.* 115, 112–124. doi:10.1016/j.ijpsycho.2017.03.002

Arbet, J., McGue, M., and Basu, S. (2020). A robust and unified framework for estimating heritability in twin studies using generalized estimating equations. *Statistics Med.* 39, 3897–3913. doi:10.1002/sim.8564

Baker, D. H., Vilidaite, G., Lygo, F. A., Smith, A. K., Flack, T. R., Gouws, A. D., et al. (2021). Power contours: optimising sample size and precision in experimental psychology and human neuroscience. *Psychol. Methods* 26, 295–314. doi:10.1037/met0000337

Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects models using lme4. J. Stat. Softw. 67, 1–48. doi:10.18637/jss.v067.i01

Bates, T. C., Maes, H., and Neale, M. C. (2019). Umx: twin and path-based structural equation modeling in R. *Twin Res. Hum. Genet.* 22, 27-41. doi:10.1017/tbg.2019.2

Blokland, G. A. M., McMahon, K. L., Thompson, P. M., Martin, N. G., de Zubicaray, G. I., and Wright, M. J. (2011). Heritability of working memory brain activation. *J. Neurosci.* 31, 10882–10890. doi:10.1523/JNEUROSCI.5334-10.2011

Bürkner, P. C. (2017). Brms: an R package for bayesian multilevel models using stan. J. Stat. Softw. 80, 1–28. doi:10.18637/jss.v080.i01

Chen, G., Pine, D. S., Brotman, M. A., Smith, A. R., Cox, R. W., and Haller, S. P. (2021). Trial and error: a hierarchical modeling approach to test-retest reliability. *NeuroImage* 245, 118647. doi:10.1016/j.neuroimage.2021.118647

Chen, G., Pine, D. S., Brotman, M. A., Smith, A. R., Cox, R. W., Taylor, P. A., et al. (2022). Hyperbolic trade-off: the importance of balancing trial and subject sample sizes in neuroimaging. *NeuroImage* 247, 118786. doi:10.1016/j.neuroimage.2021. 118786

the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study (https://abcdstudy. org), held in the NIMH Data Archive (NDA).

## Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

### Generative AI statement

The author(s) declare that no Generative AI was used in the creation of this manuscript.

### Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

### Supplementary material

The Supplementary Material for this article can be found online at: https://www.frontiersin.org/articles/10.3389/fgene.2025.1522729/ full#supplementary-material

Chen, X., Formisano, E., Blokland, G. A. M., Strike, L. T., McMahon, K. L., de Zubicaray, G. I., et al. (2019). Accelerated estimation and permutation inference for ACE modeling. *Hum. Brain Mapp.* 40, 3488–3507. doi:10.1002/hbm.24611

Chung, Y., Rabe-Hesketh, S., Dorie, V., Gelman, A., and Liu, J. (2013). A nondegenerate penalized likelihood estimator for variance parameters in multilevel models. *Psychometrika* 78, 685–709. doi:10.1007/s11336-013-9328-2

Downes, S. M., and Matthews, L. (2020). Heritability in *The stanford encyclopedia of philosophy*. Editor E. N. Zalta spring 2020 edn (Metaphysics Research Lab, Stanford University).

Downes, S. M., and Turkheimer, E. (2022). An early history of the heritability coefficient applied to humans (1918–1960). *Biol. Theory* 17, 126–137. doi:10.1007/s13752-021-00392-9

Eaves, L. J., Long, J., and Heath, A. C. (1986). A theory of developmental change in quantitative phenotypes applied to cognitive development. *Behav. Genet.* 16, 143–162. doi:10.1007/BF01065484

Falconer, D. S. D. S., and MacKay, T. F. C. (1996). Introduction to quantitative genetics. Harlow: Prentice Hall.

Fan, J., Wu, Y., Fossella, J. A., and Posner, M. I. (2001). Assessing the heritability of attentional networks. *BMC Neurosci.* 2, 14. doi:10.1186/1471-2202-2-14

Fisher, R. A. (1954). "Statistical methods for research workers," in *Biological monographs and manuals*. 12th ed., rev ed. Edinburgh: Oliver & Boyd.

Fiske, D. W., and Rice, L. (1955). Intra-individual response variability. *Psychol. Bull.* 52, 217–250. doi:10.1037/h0045276

Germine, L., Russell, R., Bronstad, P. M., Blokland, G. A., Smoller, J. W., Kwok, H., et al. (2015). Individual aesthetic preferences for faces are shaped mostly by environments, not genes. *Curr. Biol. CB* 25, 2684–2689. doi:10.1016/j.cub.2015.08.048

Guo, G., and Wang, J. (2002). The mixed or multilevel model for behavior genetic analysis. *Behav. Genet.* 32, 37–49. doi:10.1023/a:1014455812027

Gustavson, D. E., Nayak, S., Coleman, P. L., Iversen, J. R., Lense, M. D., Gordon, R. L., et al. (2023). Heritability of childhood music engagement and associations with language and executive function: insights from the adolescent brain cognitive development (abcd) study. *Behav. Genet.* 53, 189–207. doi:10.1007/s10519-023-10135-0

Haines, N., Kvam, P. D., Irving, L. H., Smith, C., Beauchaine, T. P., Pitt, M. A., et al. (2020). Learning from the reliability paradox: how theoretically informed generative models can advance the social. *Behav. Brain Sci.* doi:10.31234/osf.io/xr7y3

Harper, J., Malone, S. M., and Iacono, W. G. (2019). Target-related parietal P3 and medial frontal theta index the genetic risk for problematic substance use. *Psychophysiology* 56, e13383. doi:10.1111/psyp.13383

Holst, K. K., Scheike, T. H., and Hjelmborg, J. B. (2016). The liability threshold model for censored twin data. *Comput. Statistics and Data Analysis* 93, 324–335. doi:10.1016/j. csda.2015.01.014

Hung, I. T., Ganiban, J. M., and Saudino, K. J. (2023). Using the flanker task to examine genetic and environmental contributions in inhibitory control across the preschool period. *Behav. Genet.* 53, 132–142. doi:10.1007/s10519-022-10129-4

Hunter, M. D. (2021). Multilevel modeling in classical twin and modern molecular behavior genetics. *Behav. Genet.* 51, 301-318. doi:10.1007/s10519-021-10045-z

Iacono, W. G., Heath, A. C., Hewitt, J. K., Neale, M. C., Banich, M. T., Luciana, M. M., et al. (2018). The utility of twins in developmental cognitive neuroscience research: how twins strengthen the ABCD research design. *Dev. Cogn. Neurosci.* 32, 30–42. doi:10. 1016/j.dcn.2017.09.001

Kastrati, G., Rosén, J., Thompson, W. H., Chen, X., Larsson, H., Nichols, T. E., et al. (2022). Genetic influence on nociceptive processing in the human brain—a twin study. *Cereb. Cortex* 32, 266–274. doi:10.1093/cercor/bhab206

Maes, H. H. (2005). "Ace model," in *Encyclopedia of statistics in behavioral science*. Editors B. S. Everitt and D. C. Howell (John Wiley and Sons, Ltd), 603–605.

Martin, N. G., Eaves, L. J., Kearsey, M. J., and Davies, P. (1978). The power of the classical twin study. *Heredity* 40, 97–116. doi:10.1038/hdy.1978.10

Matthews, S. C., Simmons, A. N., Strigo, I., Jang, K., Stein, M. B., and Paulus, M. P. (2007). Heritability of anterior cingulate response to conflict: an fMRI study in female twins. *NeuroImage* 38, 223–227. doi:10.1016/j.neuroimage.2007.07.015

McArdle, J. J., and Prescott, C. A. (2005). Mixed-effects variance components models for biometric family analyses. *Behav. Genet.* 35, 631–652. doi:10.1007/s10519-005-2868-1

Mehta, P. D., and Neale, M. C. (2005). People are variables too: multilevel structural equations modeling. *Psychol. Methods* 10, 259–284. doi:10.1037/1082-989X.10.3.259

Neale, M. C., Hunter, M. D., Pritikin, J. N., Zahery, M., Brick, T. R., Kirkpatrick, R. M., et al. (2016). OpenMx 2.0: extended structural equation and statistical modeling. *Psychometrika* 81, 535–549. doi:10.1007/s11336-014-9435-8

Pingault, J. B., Allegrini, A. G., Odigie, T., Frach, L., Baldwin, J. R., Rijsdijk, F., et al. (2022). Research Review: how to interpret associations between polygenic scores, environmental risks, and phenotypes. *J. Child Psychol. Psychiatry* 63, 1125–1139. doi:10.1111/jcpp.13607

Polderman, T. J. C., Benyamin, B., de Leeuw, C. A., Sullivan, P. F., van Bochoven, A., Visscher, P. M., et al. (2015). Meta-analysis of the heritability of human traits based on fifty years of twin studies. *Nat. Genet.* 47, 702–709. doi:10.1038/ng.3285

Polk, T. A., Park, J., Smith, M. R., and Park, D. C. (2007). Nature versus nurture in ventral visual cortex: a functional magnetic resonance imaging study of twins. *J. Neurosci.* 27, 13921–13925. doi:10.1523/JNEUROSCI.4001-07.2007

Rea-Sandin, G., Clifford, S., Doane, L. D., Davis, M. C., Grimm, K. J., Russell, M. T., et al. (2023). Genetic and environmental links between executive functioning and effortful control in middle childhood. *J. Exp. Psychol. General* 152, 780–793. doi:10. 1037/xge0001298

Rijsdijk, F. V., and Sham, P. C. (2002). Analytic approaches to twin data using structural equation models. *Briefings Bioinforma*. 3, 119–133. doi:10.1093/bib/3.2.119

Robette, N., Génin, E., and Clerget-Darpoux, F. (2022). Heritability: what's the point? What is it not for? A human genetics perspective. *Genetica* 150, 199–208. doi:10.1007/s10709-022-00149-7

Rouder, J.N., Kumar, A., and Haaf, J. M. (2023). Why many studies of individual differences with inhibition tasks may not localize correlations. *Psychonomic Bulletin Review* 30, 2049–2066. doi:10.3758/s13423-023-02293-3

Rouder, J. N., and Haaf, J. M. (2019). A psychometrics of individual differences in experimental tasks. *Psychonomic Bull. and Rev.* 26, 452–467. doi:10.3758/s13423-018-1558-y

Routledge, K. M., Williams, L. M., Harris, A. W. F., Schofield, P. R., Clark, C. R., and Gatt, J. M. (2018). Genetic correlations between wellbeing, depression and anxiety symptoms and behavioral responses to the emotional faces task in healthy twins. *Psychiatry Res.* 264, 385–393. doi:10.1016/j.psychres.2018.03.042

Schachar, R. J., Forget-Dubois, N., Dionne, G., Boivin, M., and Robaey, P. (2010). Heritability of response inhibition in children. *J. Int. Neuropsychological Soc.* 17, 238–247. doi:10.1017/S1355617710001463

Schwabe, I., Gu, Z., Tijmstra, J., Hatemi, P., and Pohl, S. (2019). Psychometric modelling of longitudinal genetically informative twin data. *Front. Genet.* 10, 837. doi:10.3389/fgene.2019.00837

Sham, P. C., Purcell, S. M., Cherny, S. S., Neale, M. C., and Neale, B. M. (2020). Statistical power and the classical twin design. *Twin Res. Hum. Genet.* 23, 87–89. doi:10. 1017/thg.2020.46

Smith, D. M., Loughnan, R., Friedman, N. P., Parekh, P., Frei, O., Thompson, W. K., et al. (2023). Heritability estimation of cognitive phenotypes in the ABCD Study<sup>®</sup> using mixed models. *Behav. Genet.* 53, 169–188. doi:10.1007/s10519-023-10141-2

Smolker, H. R., Wang, K., Luciana, M., Bjork, J. M., Gonzalez, R., Barch, D. M., et al. (2022). The Emotional Word-Emotional Face Stroop task in the ABCD study: psychometric validation and associations with measures of cognition and psychopathology. *Dev. Cogn. Neurosci.* 53, 101054. doi:10.1016/j.dcn.2021.101054

Spearman, C. (1904). The proof and measurement of association between two things. *Am. J. Psychol.* 15, 72–101. doi:10.2307/1412159

Stins, J. F., de Sonneville, L. M. J., Groot, A. S., Polderman, T. C., van Baal, C. G. C. M., and Boomsma, D. I. (2005). Heritability of selective attention and working memory in preschoolers. *Behav. Genet.* 35, 407–416. doi:10.1007/s10519-004-3875-3

van den Berg, S. M., Glas, C. A. W., and Boomsma, D. I. (2007). Variance decomposition using an IRT measurement model. *Behav. Genet.* 37, 604–616. doi:10.1007/s10519-007-9156-1

van den Oord, E. J. (2001). Estimating effects of latent and measured genotypes in multilevel models. *Stat. Methods Med. Res.* 10, 393–407. doi:10.1177/096228020101000603

van der Meulen, M., Steinbeis, N., Achterberg, M., van Ijzendoorn, M. H., and Crone, E. A. (2018). Heritability of neural reactions to social exclusion and prosocial compensation in middle childhood. *Dev. Cogn. Neurosci.* 34, 42–52. doi:10.1016/j. dcn.2018.05.010

van Drunen, L., Dobbelaar, S., van der Cruijsen, R., van der Meulen, M., Achterberg, M., Wierenga, L. M., et al. (2021). The nature of the self: neural analyses and heritability estimates of self-evaluations in middle childhood. *Hum. Brain Mapp.* 42, 5609–5625. doi:10.1002/hbm.25641

Vellani, V., Garrett, N., Gaule, A., Patil, K. R., and Sharot, T. (2022). Quantifying the heritability of belief formation. Sci. Rep. 12, 11833. doi:10.1038/s41598-022-15492-0

Viktorsson, C., Lindskog, M., Li, D., Tammimies, K., Taylor, M. J., Ronald, A., et al. (2022). Infants' sense of approximate numerosity: heritability and link to other concurrent traits. *Dev. Sci. n/a* 26, e13347. doi:10.1111/desc.13347

Wang, L., Wang, Y., Xu, Q., Liu, D., Ji, H., Yu, Y., et al. (2020). Heritability of reflexive social attention triggered by eye gaze and walking direction: common and unique genetic underpinnings. *Psychol. Med.* 50, 475–483. doi:10.1017/S003329171900031X

Yeom, D., Tan, Y. T., Haslam, N., Mosing, M. A., Yap, V. M. Z., Fraser, T., et al. (2022). Genetic factors and shared environment contribute equally to objective singing ability. *iScience* 25, 104360. doi:10.1016/j.isci.2022.104360