



Published in final edited form as:

*Collabra Psychol.* 2023 ; 9(1): . doi:10.1525/collabra.77837.

## Measuring CHAOS? Evaluating the short-form Confusion, Hubbub And Order Scale

Sally A. Larsen<sup>\*,1</sup>, Kathryn Asbury<sup>2</sup>, William L. Coventry<sup>3</sup>, Sara A. Hart<sup>4,5</sup>, Callie W. Little<sup>5</sup>, Stephen A. Petrill<sup>6</sup>

<sup>1</sup>School of Education, University of New England, Australia

<sup>2</sup>Department of Education, University of York, UK

<sup>3</sup>School of Psychology, University of New England, Australia

<sup>4</sup>Department of Psychology, Florida State University, USA

<sup>5</sup>Florida Center for Reading Research, Florida State University, USA

<sup>6</sup>Department of Psychology, Ohio State University, USA

### Abstract

The Confusion, Hubbub and Order Scale (CHAOS) – short form – is a survey tool intended to capture information about home environments. It is widely used in studies of child and adolescent development and psychopathology, particularly twin studies. The original long form of the scale comprised 15 items and was validated in a sample of infants in the 1980s. The short form of the scale was developed in the late 1990s and contains six items, including four from the original scale, and two new items. This short form has not been validated and is the focus of this study. We use five samples drawn from twin studies in Australia, the UK, and the USA, and examine measurement invariance of the CHAOS short-form. We first compare alternate confirmatory factor models for each group; we next test between-group configural, metric and scalar invariance; finally, we examine predictive validity of the scale under different conditions. We find evidence that a two-factor configuration of the six items is more appropriate than the commonly used one-factor model. Second, we find measurement non-invariance across groups at the metric invariance step, with items performing differently depending on the sample. We also find inconsistent results in tests of predictive validity using family-level socioeconomic status and academic achievement as criterion variables. The results caution the continued use of the short-form CHAOS in its current form and recommend future revisions and development of the scale for use in developmental research.

\*Corresponding Author: Sally A. Larsen, School of Education, University of New England, Armidale, NSW, 2351. slarsen3@une.edu.au.

#### Contributions

Conception and design – S.A.L. and C.W.L.

Acquisition of Data – K.A., W.L.C., S.A.H. and S.A.P.

Analysis and interpretation of data – S.A.L. and W.L.C.

Drafted and revised the article – S.A.L.

Final article approval – S.A.L., K.A., W.L.C., S.A.H., C.W.L. and S.A.P.

Competing Interests

None

## Keywords

Confusion Hubbub and Order Scale; home environment; factor analysis; measurement invariance; twin studies

---

The effects of home environments on childhood functioning and development has been a topic of research interest for decades (Bradley, 2015; Evans, 2006). Bronfenbrenner's widely-known bioecological model of human development defines the home environment as a key context for proximal processes that influence childhood development (Bronfenbrenner, 1986; Bronfenbrenner & Ceci, 1994; Bronfenbrenner & Morris, 2006). In particular, Bronfenbrenner argued that stability in home environments was particularly important for development (Bronfenbrenner & Evans, 2000), and disruptions to routine family life were centrally important in poor childhood psychological functioning (Bronfenbrenner & Ceci, 1994).

Identifying the features of stable and consistent family home environments, and examining the effects of variability in home contexts, has therefore been important for testing the propositions of the bioecological theory of development. To achieve this aim, however, home environmental features must be recorded or measured in some way. The purpose of the current study, therefore, was to examine the measurement properties of a widely-used short-form scale, which was developed with the intention of capturing variability in home environments, the Confusion, Hubbub And Order Scale (CHAOS; Matheny et al., 1995). We begin with some background on the measurement of home environments before providing a brief history of the CHAOS measure.

## Capturing the variability in home environments

In the mid-20<sup>th</sup> century researchers recorded information about differences between home environments via in-person observations (e.g. Wilson & Matheny, 1983). Research assistants spent many thousands of hours attempting to unobtrusively observe aspects of home environments, including interactions between parents and children, the number of visitors coming and going, ambient noise levels (both internal and external to the home), observable routines established by parents, among other features (Evans, 2006). These efforts to measure the quality of home environments were expensive, and limited in that a finite number of households in a geographically constrained area could be visited by research assistants within a given timeframe. In the 1980s, therefore, measurement of household environments began to shift from a reliance on observations, recorded as both qualitative information and observer ratings on quantitative scales, to self-report scales, where parents were asked to rate aspects of their homes according to questionnaire items.

The CHAOS scale was one of several self-report instruments established during the latter decades of the 20<sup>th</sup> century (the HOME scale is another widely used instrument; c.f. Bradley, 2015 for a review). Two versions of the CHAOS scale have been used in research since its inception: the original 15-item version proposed by Matheny et al. (1995), and a short-form 6-item version. The scale has been used extensively in research: there are over 500 citations to date of the paper reporting the psychometric properties of the original scale

(Matheny et al.), and a library database search identified 305 articles and 62 dissertations referencing the name of the scale (June, 2022). It is difficult to identify how many of these research articles use the long form and how many use the short form of the scale given that both have the same title. The latter, short-form is the focus of this study, however some background on the original long-form is relevant here.

The 15-item CHAOS measure was developed in a sample of over 400 families of twins participating in the Louisville Twin Study during the 1980s. The original scale (reproduced in Figure 1.) comprised true/false scored items, half reverse scored, which were summed to produce an overall measure of the (in)stability of the household environment. Conceptually the scale captured aspects of household confusion and disorder, including high levels of noise, clutter, disorganization and “frenetic activities” (p.432). Matheny et al. (1995) validated the 15-item scale using a subsample of 123 mothers of infants ranging in age from 6 to 30 months reporting a reliability coefficient of  $\alpha = .79$ . A further subsample of 42 mothers completed the questionnaire at a 12-month interval with a test-retest correlation of  $r=.74$ . Matheny et al. noted that the 15-item CHAOS scale accounted for a unique proportion of systematic within-home differences that could not be attributed to parent education or SES measures. Nonetheless, direct observations of the home environment were not interchangeable with the CHAOS scale: there was a significant, but not high, degree of overlap between the two measures in the study ( $R^2 = .39$ ). The 15-item scale was further validated in two different samples of preschool ( $n = 106$ ) and school-aged children ( $n = 676$ ; Dumas et al., 2005), demonstrating overlap with, but distinction from, measures of socioeconomic status, and good internal consistency reliability ( $\alpha = .83 / .81$ ).

The short 6-item form of the CHAOS measure was first used in the late 1990s by studies including the Twins Early Development Study (TEDS; e.g. Asbury et al., 2003) and the Western Reserve Reading and Math Project (WRRMP; e.g. Hart et al., 2006). The short-form of the scale is reproduced in Figure 2. This form consists of 4 items from the original Matheny et al. (1995) 15-item scale, plus the addition of two items not appearing in the original scale: Item 1 *The children have a regular bedtime routine* and item 5 *There is usually a television turned on somewhere in our home*. There is no published information on why these 6 items were chosen for inclusion in the short-form CHAOS and no evaluations of whether the short form captures the full range of the original intended construct (e.g. Smith et al., 2000). In published articles using the short-form, construct validity evidence is universally attributed to the article which reports the validity of the long-form, 15-item scale (Matheny et al., 1995), and reliability information is usually reported as Cronbach’s coefficient alpha, with estimates ranging from  $\alpha = 0.52 - 0.68$ .

Notwithstanding the lack of published evidence that the short-form CHAOS scale reliably measures the same construct as that proposed by the long-form, many studies have used the 6-item scale to examine links between the home environment and childhood functioning. Table S1 in the supplementary material shows details of 21 papers that we could identify published between 2003 and 2019 that used the short-form CHAOS. Studies have examined associations between CHAOS and cognitive development (Petrill et al., 2004; Pike et al., 2006), reading skills (Johnson et al., 2008), language development (Asbury et al., 2005), and behavioural problems (Coldwell et al., 2006; Deater-Deckard et al., 2009; Laurent et

al., 2014; Peviani et al., 2019), consistently finding that higher ratings of household CHAOS are associated with worse functioning or development. The short scale has also been used in studies examining social determinants of health and wellbeing (Ganasegeran et al., 2017; Suku et al., 2019). Many of the studies that have collected data on the CHAOS short-form have been twin studies which examine home environments within the behaviour genetics theoretical framework: home environments are considered an aspect of *shared environments*, i.e. environmental features which serve to make twins more similar to one another (Plomin et al., 2013), although it is acknowledged that twins can perceive the same objective environment differently (Hanscombe et al., 2010). Several studies have examined the extent to which home environments mediate or moderate genetic influences on childhood outcomes (Asbury et al., 2003; Gould et al., 2018; Harlaar et al., 2005; Hart et al., 2007; Petrill et al., 2004), and one study attempted unsuccessfully to identify genetic influences on reports of CHAOS using a genome-wide association design (Butcher & Plomin, 2008).

The short-form scale has also undergone several additional transformations, including translations into languages other than English (e.g. Deater-Deckard et al., 2019; Ganasegeran et al., 2017), and versions where children or adolescents themselves rate their home environments on a three-point likert response scale. Using this adolescent self-report data, studies have examined links between CHAOS and academic achievement, behavioural functioning (Hanscombe et al., 2010, 2011; Kim-Spoon et al., 2017), and brain activity in functional MRI studies (Lauharatanahirun et al., 2018). An even shorter version of the short-form, comprising only five items, is also evident in the literature. This five-item version is used in the Parenting across Cultures Study, which recruited 511 urban families in six low-to-middle income countries (China, Kenya, the Philippines, Thailand, Colombia and Jordan; Deater-Deckard et al., 2019). In this version, item 5 *There is usually a television turned on somewhere in our home*, was omitted because of the possibility that families in low income countries do not own televisions. There is, as yet, no published scale evaluation information indicating this item performs badly. Despite this additional variation, the studies emerging from the Parenting Across Cultures project report the original Matheny et al. (1995) article as evidence of the reliability of the scale for capturing “a harsh and unpredictable environment” (Chang et al., 2019a, p.4; Chang et al., 2019b; Deater-Deckard et al., 2019).

Given the widespread use of the short-form CHAOS, and its attractiveness in terms of minimal time commitment of respondents in multivariate surveys, the lack of reliability and construct validity information for the scale is of concern. The present study thus aims to examine the measurement properties of the short-form CHAOS scale. We use several approaches and five datasets to adjudicate whether the scale is valid and reliable for measuring the quality of home environments. In defining validity, we take the position of Borsboom et al. (2004) who argued simply that “a test is valid if it measures what it purports to measure” (p.1061). Furthermore, we define reliability as “an index of measurement precision” (p. 1070) that can be evaluated within a scale (i.e. how well do items measure the same construct) and across measurement occasions (i.e. between samples or within samples over time). In this study we therefore: 1) examine the factor structure of the six items, 2) evaluate whether the measure is invariant across groups, and 3) examine the predictive validity of the scale using a measure of socioeconomic status and childhood

academic achievement as criterion variables. In this way we collate evidence of the validity and reliability of the scale as an adequate measure of the quality of home environments in different populations.

## Measurement invariance

Combining multiple survey items into a single composite score is very common practice in social science research. Creating a sum or average from several items, however, assumes that the scale in question captures one underlying factor (McNeish & Wolf, 2020). In cases where the construct of interest has been shown to capture a single factor, using composites is a defensible strategy (Widaman & Revelle, 2022), however, there is minimal documentation regarding the most appropriate factor structure of the short-form CHAOS (except in Johnson et al., 2008). Furthermore, use of a scale in different populations also assumes that the measure captures the same latent construct regardless of context (Millsap & Olivera-Aguilar, 2012). Any differences in the means or variances of the observed items is assumed to be related to differences in the latent construct itself, rather than differences between populations that are not related to household order and routine. In this study we test the assumptions that a) a single factor underlies the six items in the short-form CHAOS, b) the factor structure is the same across samples, and c) differences on the observed variables are caused by differences on the latent construct, and are not due to unobserved, external differences between the populations of interest.

We can begin to test these assumptions using confirmatory factor analyses and a measurement invariance procedure. We follow the typical procedure for testing measurement invariance recommended by methodologists (e.g. Millsap & Olivera-Aguilar, 2012; van de Schoot et al., 2012): namely, the same confirmatory factor model is first estimated in each group, then increasingly restrictive equality conditions are introduced for different sets of parameters. If measurement invariance holds across samples, we can be confident that comparing the results of studies using the CHAOS scale in different contexts is valid and informative. On the other hand, if the analyses indicate that the measure is non-invariant, response patterns on the observed items could be systematically influenced by unobserved differences between populations, for example interpretive differences for specific items, rather than by differences in the latent domain of interest (i.e. household order and routine).

## This study

Our hypotheses for the study are informed by, 1) the consistently low reliability reported in studies using the short-form CHAOS ( $\alpha = .52 - .68$ ), 2) evidence from one study that the six items are better represented by two factors rather than one (Johnson et al., 2008), and 3) preliminary evidence generated by an exploratory factor analysis indicating a two-factor solution (see below). Given this information, we hypothesised that a two-factor dimensional structure will better fit the data on the CHAOS items in all samples. Preliminary analyses also inform our hypothesis that the measure will be non-invariant across the five samples: that is, we do not expect the six items to behave similarly in all samples, nor do we expect the factor structure to be repeatable across samples. Finally based on the research findings described above we predict that higher CHAOS will be negatively

associated with both family socioeconomic status and academic achievement, however in the case where a two-factor model is most appropriate, it is not clear whether one or both factors will be significantly associated with each criterion variable. Johnson et al. (2008) demonstrated that only one of two factors (*household order* but not *noise*) was associated with several measures of childhood literacy, however whether the factor structure identified in this previous analysis holds across all samples will only become evident after the initial invariance testing across all five samples. Preliminary hypotheses and an overview of the study were preregistered at the Open Science Framework (<https://osf.io/akmf4>). Subsequent to preregistration we gained access to an additional dataset not noted in the preregistration (the Project KIDS data). Ethical approval for this study was obtained from the first author's institution (Human Research Ethics Committee Approval# HE22-093).

## Methods

Secondary data for the project was sourced from: three studies located in the US, the Western Reserve Reading and Math Project (WRRMP; Hart et al., 2007; Petrill et al., 2006), the Florida Twin Project on Reading, Behavior and Environment (FTP-RBE; Taylor et al., 2019), and Project KIDS (Kids and Individual Differences in Schools; van Dijk et al., 2022); one study located in the UK, the Twins Early Development Study (TEDS; Oliver & Plomin, 2007; Rimfield et al., 2019); and one study located in Australia, the Academic Development Study of Australian Twins (ADSAT; Larsen et al., 2020). These studies were selected because all collected parent reports on the English language short-form CHAOS using a five-point likert response scale (see Figure 2.), and the children of interest were aged between 3 (earliest wave of TEDS) to 12 years (upper age of FTP-R, ADSAT and Project KIDS wave 1 samples). Descriptive statistics of participants in all samples and data collection waves are in Table 1.

## Samples and Measures

*The Academic Development Study of Australian Twins (ADSAT)* recruited a national sample of 2762 families of Australian school aged twins between 2012 and 2017 (Larsen et al., 2020). The design of study recruitment was partly prospective and partly retrospective. For the current investigation we selected only families who were recruited to the study and had completed the CHAOS measure when their twins were in Grades 3, 4 or 5 ( $n=1294$ ; age 8 to 11 years). This age group was selected in an attempt to align the ages of participating children as closely as possible across samples. For an initial exploratory factor analyses we used an additional sub-sample of 596 families participating in the ADSAT who also completed the CHAOS form on enrolment into the study. Parents completed the CHAOS measure only once.

Academic achievement was measured by standardized scores on reading comprehension and mathematics tests undertaken by children as part of the National Assessment Program: Literacy and Numeracy (NAPLAN; ACARA, 2017). For this study we used scores on the Grade 3 assessments to align with when the CHAOS items were collected. Socioeconomic status in this dataset is a factor score comprising the highest educational attainment of both

parents, the occupational prestige ranking of both parents and an indicator of neighbourhood socioeconomic advantage (see Larsen et al., 2020 for details).

The *Florida Twin Project on Reading, Behavior and Environment (FTP-RBE)* is a subsample of the 2753 twin pairs recorded in the Florida State Twin Registry (FSTR; Taylor et al., 2019). Beginning in 2012, a subsample of families with twins enrolled in the FSTR were invited to enrol in the FTP-RBE, which involved completing a questionnaire, containing in part the CHAOS items, every other year over six years (i.e., three waves of questionnaire assessment). The mean age of twins for the first wave of the questionnaire data collection was 11.16 years. In total, 568 families (72% of the invited participants) provided data on the CHAOS at wave 1, reducing to 447 at wave 2 and 313 at wave 3. Academic achievement was measured by scores on the Florida Comprehensive Assessment Tests (FCAT) reading subtest, undertaken by students in the 2011–12 and 2012–13 school years. The FCAT test is a standardized assessment of reading, completed by students at the end of grades 3 to 11. FCAT data were provided by Florida's Progress Monitoring and Reporting Network (PMRN). Socioeconomic status is a factor score generated using five observed variables: estimated family income, both parents' highest educational attainment, and both parents' occupational prestige.

*Project KIDS* is a repository of data collected in nine randomized control trials of reading interventions undertaken in the US between 2005 and 2011 (see van Dijk et al., 2022). Data on the CHAOS short form was collected in 2013 from a sub-sample of 442 families of singleton children who had participated in at least one trial. Data on school achievement was collected in the same parent survey. For both English Language Arts and Math, parents reported their children's achievement on a 5-point rating scale, ranging from *A/Excellent* (1) to *F/Fail* (5). Achievement variables were reverse coded before analysis so that higher ratings indicated better achievement, similar to other achievement tests used in this study. Socioeconomic status observed variables and factor score estimation was identical to that in the FTP-RBE study described above.

*Western Reserve Reading and Math Project (WRRMP)* is a longitudinal cohort-sequential study which recruited families of twins, primarily in the state of Ohio, USA, beginning in 2002. Twins were in kindergarten or first grade on recruitment (mean age = 6.09 years) and were followed up on measures of literacy and CHAOS approximately annually over seven waves of data collection. Across all waves, 794 families provided at least some data to the project. The short-form CHAOS was collected at each wave of the study, with 580 families answering the items in wave 1, reducing gradually to 246 families responding by wave 7.

We selected five assessments of academic skills in both reading and math domains collected across all waves of the WRRMP. These included, a) two assessments of reading comprehension, the PIAT-R/NU (Dunn & Markwardt, 1998) and the WRMT-R passage comprehension assessments (Woodcock, 1987), and b) three assessments of math subdomain skills, the Woodcock-Johnson calculation, applied problems, and quantitative concepts tests (Woodcock, 1987). All children who were able to be followed up at each wave provided data on these assessments. For the WRRMP study we used a proxy of socioeconomic status using

variables that were available in the dataset: an average of both parents' highest educational attainment.

The *Twins Early Development Study (TEDS)* recruited a nationally representative sample of 13,732 families of infant twin pairs in the United Kingdom from 1994–1996 (Oliver & Plomin, 2007). For this study we use wave 3 and 4 of data collection, when twins were aged 3 and 4 years, respectively. Parents responded to the short-form CHAOS in both waves, with 6009 parents providing data in wave 3, and 8014 in wave 4. Later collections on the CHAOS measure used a 3-point likert response scale and/or asked twins themselves to respond, rather than parents. We omit these waves and focus on the CHAOS measure obtained in the same manner as that for the other data collections in this study. Academic achievement was assessed when twins were aged approximately 7 years. All students in the UK undertake National Curriculum assessments in core subjects. Standardized assessment results for English and Mathematics were sourced from government data collections. We use the socioeconomic status variable available in the TEDS dataset, a composite variable generated from five variables: occupational prestige of both parents, highest educational levels of both parents, and mother's age at the birth of the first child.

We note that four of the five studies included in this project were studies of child and adolescent twins. For each dataset, the CHAOS items and SES variables were collected at the family level (i.e. one set of responses by family for each wave in each study), therefore we did not need to account for the nested nature of data collected on twin pairs. Academic achievement variables were collected for each twin separately, however, so in instances where we use achievement as criterion variables, we selected one twin at random from each pair. We do not report results for the second randomly-selected twin, but findings were no different.

### Analysis plan

In this study we aimed to test the factor structure, measurement invariance and predictive validity of the short-form CHAOS using five samples collected in different contexts. We first wanted to test whether the usual approach to using the six items – i.e. combining them into a single mean or sum score – is the optimal approach to the use of the scale. Only one study to date has reported an exploratory factor analysis of the items (Johnson et al., 2008). Using the WRRMP Wave 1 data this study demonstrated a two-factor solution in an exploratory factor analysis (EFA) with one factor comprising items 1, 4 and 5 (termed “household order and routine”), the second factor comprising items 2, 3 and 6 (“quietness of the household”; Johnson et al., p. 5). The two factors correlated at  $r = .33$ . The proportion of variance explained by the two-factor solution and the factor loadings of the items were not reported. Items were subsequently summed within each factor for further analyses. Given this study is the only one to date to examine the factor structure of the short-form CHAOS, the first step in the analysis for the current study was an EFA using a subsample of participants in the ADSAT ( $n=596$ ). Specifically, a principal components analysis using full information maximum likelihood estimation was undertaken. Due to the results in Johnson et al. we expected that a two-factor solution would fit the data better than a one-factor



model. Therefore, we examined eigenvalues, compared the proportion of variance explained by one- and two-factor solutions, and examined item-factor loadings.

To further examine whether a one- or two-factor structure of the six items was best supported by all datasets, we next ran confirmatory factor analyses (CFA) separately for each sample. For CFAs comprising two factors we allowed factors to correlate, but did not allow any cross-loadings of items, nor any residual covariances. We examined model fit statistics, and compared nested models to identify the best solution in each sample. We predicted that two-factor models would be a better fit to the data than one-factor models for all samples, however we made no specific predictions about whether the configuration of items reported by Johnson et al. (2008) would be the best fit in each sample.

Next, measurement invariance was examined via multiple-group confirmatory factor models. In this step we consider each dataset a different group since each study was conducted in a different context, and three countries are represented by the five datasets, Australia, the UK and the USA. We followed the procedure suggested in several sources and tested i) configural, ii) metric, and iii) scalar invariance (e.g. Byrne, 2012; Meredith & Teresi, 2006; Millsap & Olivera-Aguilera, 2012; Putnick & Bornstein, 2016). We did not expect strict invariance (i.e. invariance of residuals) to hold across groups so planned to test this step only where scalar invariance was confirmed. Specifically, configural invariance models force the same factor structure across groups but allow item loadings, item intercepts and residuals to vary. Because we planned to first test confirmatory factor models for each group, and select the best-fitting model, we expected configural invariance to hold. Metric invariance forces equivalence of factor loadings across groups and assesses whether this restriction leads to a significant reduction in model fit. Scalar invariance tests for equivalence of item intercepts across groups retaining the equivalence of factor loadings tested in the previous step. Strict invariance retains the equivalence constraints introduced by metric and scalar invariance, and constrains item residuals to equality. If at any step model fit statistics suggested significantly poorer fit, we examined the parameters constrained by that step to identify potential sources of model misfit.

Model fit was assessed using several statistics. Given that  $\chi^2$  goodness of fit is affected by large samples or variable sample sizes in multiple group models (Byrne, 2012), we report this statistic along with several others. In particular, we examine the root mean square error of approximation (RMSEA), which ideally should fall  $< 0.08$  (Byrne, 2012). We also examine the comparative fit index (CFI), which provides an estimate of incremental fit of the model compared with a baseline model. Current advice suggests CFI values of  $> 0.95$  indicate adequate model fit (West et al., 2012).

For assessing the model fit of the nested models, such as those in each step of the measurement invariance tests, we examine the change in  $\chi^2$  relative to change in degrees of freedom ( $df$ ). Ideally the change in  $\chi^2$  for each  $df$  should have  $p > .001$ , indicating that the more restricted model is not a worse fit to the data than the less restricted model. When equating parameters across groups in measurement invariance analyses, particularly when large numbers of groups are compared, RMSEA and CFI can also be examined (Rutkowski & Svetina, 2014). A change of 0.010 (RMSEA) and  $-0.010$  (CFI) are indicative

of non-invariance between groups when parameters are constrained to equality for metric or scalar invariance tests (Cheung & Rensvold, 2002; OECD, 2010). Finally, Akaike's Information Criterion (AIC) can provide additional information about fit for non-nested models with smaller values indicating better model fit. We report and interpret AIC where appropriate (West et al., 2012).

It is important to note that interpreting change in model fit statistics to assess measurement invariance across more than two groups, as we do in this study, can generate information without clear or simple interpretations. For example, should model fit decrease significantly at any step of measurement invariance testing, with five groups in the model, it may not be clear whether one sample is driving model misfit, while others are sufficiently comparable. Notwithstanding this interpretational problem, the main aim of the study is to evaluate whether the CHAOS measure behaves similarly across contexts, therefore non-invariance of even one sample is problematic for the applicability, use, and interpretation of the scale in different contexts.

Finally, we planned to examine the predictive validity of the CHAOS measure using two criterion variables. We examined zero-order correlations between CHAOS and a socioeconomic status variable (or proxy), and academic achievement variables available in each dataset. We compared results using a) a one-factor model of CHAOS, b) a two-factor model, and c) analyses where the CHAOS items are composed as factor scores, with results when items are composed as mean scores, as is more common in the published literature.

All analyses were run in the statistical program R (R Core Team, 2020) using the psych package (Revelle, 2022) for descriptive statistics, reliability statistics, creating factor scores and exploratory factor analyses, the lavaan package (Rosseel, 2012) for confirmatory factor models and invariance testing, and ggplot2 (Wickham, 2016) for figures. Code for confirmatory factor analyses, and invariance testing is at the OSF (<https://osf.io/akmf4>). Data from FTP-RBE, Project KIDS and the WRRMP is available at LDBase repository (Hart et al., 2020). Data from the ADSAT is available on request to the first author, and data from TEDS is available on request from data managers (Kings College London, 2022, <https://www.teds.ac.uk/researchers/teds-data-access-policy>).

## Results

For each sample means, standard deviations, skew, and kurtosis of each item, and zero-order correlations between items were generated. These are reported in Tables S2–S11 in the supplementary material. We report three waves of data for multi-wave studies, except TEDS, which contains only two waves of parent-report on the CHAOS scale. Correlations between items were all positive, with some variation in the strength of correlations across the samples. Perhaps most notable were the differences in correlations between item 5. *There is usually a television on somewhere in our home*, and the remaining items. In the Project KIDS and FTP-RBE samples, correlations between this item and the remaining five were generally smaller ( $r = .17$ ) than those in the ADSAT, WRRMP and TEDS samples ( $r = .34$ ). On the other hand, the strongest correlation in all samples was between items 2. *You can't hear yourself think in our home* (reversed) and 3. *It's a real zoo in our home* ( $r = .56 - .77$ ).

Variation is also evident in item means and distributions across studies. Figures S1–S5 (supplementary materials) show item distributions for each dataset when selecting one wave from each multi-wave study. Response patterns were similar over waves within studies (i.e. for TEDS, WRRMP and FTP-RBE). The mean of item 2, *You can't hear yourself think in our home* (reversed) varied from 1.95 (Project KIDS) to 3.29 / 3.27 (TEDS sample, wave 1 / wave 2). Similarly, item 5, *There is usually a television on somewhere in our home*, also shows variable response patterns across studies, as does item 3, *It's a real zoo in our home*. Finally of note is the strong agreement with item 1, *The children have a regular bedtime routine* across all studies (i.e. most respondents selected *Somewhat true* or *Definitely true* for this question).

Given the six items in the scale are most often used as a sum or average score, we also computed coefficients alpha ( $\alpha$ ) and omega (hierarchical,  $\omega_h$ ). These are reported in Table 1 and range from  $\alpha = 0.55$ – $0.70$ , and  $\omega_h = 0.29$ – $0.63$ , indicating that internal consistency for the six-item scale is poor across samples. In particular, given that  $\omega_h$  is arguably a more appropriate indicator of reliability because it allows for different factor loadings of items (McNeish, 2018), the coefficients of  $<0.65$ , with three  $<0.50$ , suggested that a one-factor model for the six items would not be supported in confirmatory factor analyses.

### Exploratory factor analysis.

To support decisions about whether a one-factor or a two-factor solution would be most appropriate across all samples, we next undertook an exploratory factor analysis (EFA) on a subsample of participants in the ADSAT ( $n=596$ ). We selected the ADSAT data for the EFA because the first author had access to these data before obtaining permission to access the remaining four datasets. Table S12 (in the supplementary materials) shows eigenvalues and proportion of variance explained for each principal component in the EFA, estimated using the maximum likelihood method. The first two components collectively explained 57% of the variance, and both had eigenvalues  $> 1$ . Remaining components had eigenvalues  $< 1$ , and the parallel analysis plot (Figure S6 in supplementary materials) also supported a two-factor solution.

Factor loadings and communalities for a two-factor solution are reported in Table S12 (supplementary material). Interestingly the EFA in this sample suggested a different pattern of items loading on each of two factors to that indicated by the published example using the WRRMP data (Johnson et al., 2008). In the WRRMP data, items 1, 4, and 5 comprised one factor, termed '*order and routine*', and items 2, 3 and 6 comprised the second factor, labelled '*quietness of the household*'. In the ADSAT data, by contrast, items 2, 3 and 5 loaded on one factor and appeared to represent *household noise*, while item 6 cross-loaded on both factors. This cross-loading suggests the wording of item 6, *The atmosphere in our house is calm*, could be interpreted in the light of either household noise, or routine. Given this inconsistent result, we opted to test two different configurations of the six items in confirmatory factor models. The configuration identified by Johnson et al. grouped items 1, 4 and 5 (*routine*), and 2, 3 and 6 (*quietness*). The second configuration informed by the EFA described here grouped items 1, 4 and 6 (*disorder*), and items 2, 3, and 5 (*noise*).

### Confirmatory Factor Analysis by sample.

To evaluate whether the one-factor, or either of the proposed two-factor structures of the six items was consistently reproduced over the five samples, we first tested the three models separately in all samples and waves. First, and in alignment with the common usage of the scale as a sum or average of the six items, we tested a one-factor model, forcing all items to load on one latent variable with no residual correlations (Table 2, Model A. in all samples and waves). We compared this one-factor model with the two different configurations of a two-factor model, each allowing three items to load on each factor (see justification above). Because the two-factor models are not nested (the same number of parameters is estimated in both) we compare each (i.e. Models B. and C. in each panel of Table 2.) with the one-factor model. While this comparative process is imperfect given that Models B. and C. cannot be directly compared using most fit statistics, evaluating the fit of each model against the one-factor option does provide some information on which solution may be more appropriate. In addition, an examination of the AIC provides additional information about which two-factor model might be retained.

In all samples, a one factor model (A.) was a poor fit to the data according to all criteria (Table 2., first row of each panel). In all cases the RMSEA statistic did not fall within the acceptable range, and the CFI and TLI statistics were  $<.95$ . Model B tested the two-factor solution reported in Johnson et al. (2008), with factors termed '*quietness*' and '*routine*'. Change in  $\chi^2$  (*df*), RMSEA, CFI and AIC for the two-factor model compared with the one-factor model showed an improvement in fit in all samples. Nonetheless, in most cases fit statistics were poor or borderline. The exception was the WRRMP dataset (that used by Johnson et al.), which showed borderline-good model fit for this configuration of items in five of seven waves.

Model C tested the alternative configuration of the six items suggested by the EFA in the ADSAT data. This model fit the data better than the one-factor model according to all fit statistics (Table 2, model C.). AIC statistics indicated that this alternative two-factor configuration was a better fit to the data than that tested in model B for all samples and waves except for the WRRMP data. In the ADSAT, Project KIDS, wave 2 of the WRRMP, and both waves of TEDS data, fit statistics were acceptable or borderline for model C.

Nonetheless, despite the improvement relative to model A, the fit of model C in the FTP-R data and wave 3 of the WRRMP remained poor when evaluating the RMSEA, CFI and TLI against suggested cut-off criteria. Notwithstanding this problem, we retained the first wave of the FTP-R data for multiple group invariance testing because this wave had the least missingness. We also retained wave 2 of the TEDS sample (older age group and less missingness), and wave 2 of the WRRMP data (best fit for model C.), the ADSAT and Project KIDS samples.

### Measurement Invariance.

Table 3. shows fit statistics for each step of invariance testing incorporating all five samples. The configural invariance model forces the same configuration of items loading on factors across all groups with no cross-loadings or residual covariances for observed

items. Factor loadings, intercepts, variances and covariances are allowed to vary by group. Notwithstanding the poor fit of the models for some individual samples noted above, the configural invariance model (Table 3, Model 1) showed borderline acceptable fit to the data when evaluated by the RMSEA (0.077, 90% CI [0.071, 0.083]) and CFI (0.957) statistics. Next, model 2A. (Table 3) tested for metric invariance by constraining factor loadings of all items to equivalence across groups. According to the AIC and the  $\chi^2$  difference relative to degrees of freedom ( $\chi^2$  (df) = 99.10 (16),  $p < .001$ ), the fit of model 2A was significantly worse than model 1. However, the RMSEA (0.007), and CFI (0.007) indicated the fit of this model was not worse relative to model 1 (using cutoff values of 0.010 for each; Rutkowski & Svetina, 2014). Given this mixed information, we examined the factor loadings across the five samples in the configural invariance model. The loadings for item 5 *There is usually a television turned on somewhere in our home* (reversed) were notably different across samples, ranging from 0.07 (FTP-R) and 0.09 (Project KIDS) to 0.47 (TEDS), 0.37 (WRRMP) and 0.35 (ADSAT). Consequently, we released the constraint on the loading for this item, and tested a partial metric invariance model with the remaining five item-loadings constrained to equivalence.

The partial metric invariance model (2B. in Table 3) fit the data significantly better than the full metric invariance model (2B vs 2A:  $\chi^2$  (df) = 68.36 (4)), and was not a worse fit to the data than the configural invariance model (2B vs 1:  $\chi^2$  (df) = 30.74 (12),  $p = .002$ ;

RMSEA = 0.002; CFI = 0.002). We thus retained the partial metric invariance model and next tested scalar invariance by constraining all item intercepts to equivalence across the five samples. Fit statistics for scalar invariance (model 3, Table 3) show that this model was a worse fit to the data on all criteria compared with the partial metric invariance model ( $\chi^2$  (df) = 3320.65 (16),  $p < .001$ ; RMSEA = 0.092; CFI = 0.278). We therefore retained the partial metric invariance model and examined the item intercepts by group for possible reasons why scalar invariance was not supported.

Table 4. shows factor loadings and intercepts for each dataset for the retained partial metric invariance model for five samples. There is considerable variation in intercepts for some items across the five groups, after holding loadings constant for all but one item.

For example, the intercepts for item 2, *You can't hear yourself think in our home* (reverse coded), range from 1.95 in the Project KIDS sample to 3.27 in the TEDS sample; similarly, for item 3. *It's a real zoo in our home*, intercepts range from 1.74 in the Project KIDS sample to 2.66 in the TEDS sample (N.B. because these item intercepts are allowed to vary by group, the model essentially reproduces item means reported in Tables S1–S10). The loadings for most other items have a smaller range, for example, 1.81–1.94 for item 4. *We are usually able to stay on top of things*. These differences in factor loadings indicate that response patterns vary across samples, potentially for reasons which are unrelated to differences in the latent construct under consideration (i.e. the confusion, hubbub and order of the home environment).

The r-square values reported in Table 4 provide additional information about the extent to which the variance in each item is captured by the final model. Of note is the low  $R^2$  for two items. First, for item 5. *There is usually a television turned on somewhere*

*in our home*, variance explained ranged from 0.2% (FTP-RBE), to 7% (ADSAT), to 14% (TEDS). Similarly, for item 1. *The children have a regular bedtime routine*,  $R^2$  values were persistently low, with 3–5% of the variance explained by the factor model in all datasets. It is worth noting that both these items were first introduced when the short-form CHAOS was created, and did not appear in the original 15-item scale. For remaining items  $R^2$  ranged from 16 to 77%.

### Predictive validity.

To examine the predictive validity of the short-form CHAOS, we estimated correlations between two different configurations of the six items and available academic achievement variables for each dataset. We selected the same wave of data as that selected for between group measurement invariance tests described above. First, we generated a single-variable factor score using all six items, following the most common use of the scale. Secondly, we generated factor scores for two variables based on the two-factor solution with the best-fitting model. Specifically, these factors comprised three variables each and were termed *disorder* (items 1, 4 and 6), and *noise* (items 2, 3, and 5). Using the *psych* package in R (Revelle, 2022), factor scores were generated separately for each configuration of items, producing variables with  $M=0$  and  $SD=1$ . Table 5 shows correlations between factor scores and criterion variables for one- and two-factor configurations of items. For comparative purposes, we also generated composite variables for both combinations of the six items, i.e. single variable averaging across the six CHAOS items, and two variables using averages of the same three items as used in the factor score models. Table 6 shows correlations between average CHAOS scores and criterion variables.

Correlations with achievement were either negative (as expected), or negligible, varying by sample and whether a one-factor or two-factor combination of items was used. For four of the five datasets – the ADSAT, FTP-RBE, Project KIDS and TEDS – the correlations reported in Table 5 supported our prediction that higher levels of parent-reported confusion, hubbub and disorder in homes would be negatively associated with measures of academic achievement. In these four datasets, correlations between a one-factor CHAOS measure and achievement ranged between  $-.07$  to  $-.21$ . Interestingly, when the six variables were separated into two factors, only the *Noise* factor consistently correlated with achievement, while the *Disorder* factor did not. The most notable exception was the Project KIDS data where the *Disorder* factor correlated negatively with both English and Math grades ( $r = -.20 / -.17$  respectively), even though the correlation between Math grades and the one-factor CHAOS was small and not significantly different from zero ( $r = -.06$ ). By contrast, the correlations between one-factor CHAOS variable and the two factors (*Noise* and *Disorder*) were small and non-significant for the WRRMP across several high-quality measures of both reading comprehension and mathematics sub-domain skills (calculation, applied problems, and quantitative concepts).

Results were similar to those reported for the factor scores when composite variables were created by averaging across items (see Table 6). One exception was noted in the WRRMP data, where the WRMT passage comprehension measure correlated negatively with the

composite score of six items ( $r = -.16$ ) and with the *Noise* composite ( $r = -.20$ ), and three of four other reading and math measures correlated negatively with the *Noise* composite.

Correlations between SES and CHAOS also differed by sample, whether one or two scores were used, and method of creating variables (factors or averages). When average scores were used (Table 6), SES correlated negatively with the one-variable CHAOS across all datasets ( $r = -.15$  to  $-.32$ ). Similar to the academic achievement variables, when CHAOS was separated into the two composites, SES correlated more consistently with the *Noise* composite and not the *Order* composite. When factor scores were used to generate the CHAOS variables (Table 5), patterns of correlations with SES were similar, but notably smaller in all datasets, and not significantly different from zero in the WRRMP sample. In the Project KIDS sample the *Noise* factor alone correlated with CHAOS ( $r = -.23$ ,  $p < .001$ ) even though the composite using all six CHAOS items did not ( $r = -.09$ ,  $p = .06$ ).

## Discussion

The central aim of this study was to examine the measurement properties of the short form of the Confusion, Hubbub and Order Scale (CHAOS) with a goal to provide some recommendations on the use of the scale – both in pre-existing datasets and new research. On the whole, our results indicate that the six items in the short-form CHAOS are not reliable and valid *enough* to capture variability in the quality of home environments across different contexts, age ranges of child study participants, and time. Ideally, reduction of a long-form to a short-form scale should be accompanied by evidence that the short form itself is a) reliable, b) valid, and c) captures the breadth of the construct indicated by a long-form (Clarke & Watson, 1995, 2019; Smith et al., 2000). Because these steps were not documented for the short-form CHAOS, this study thus provides some of the information necessary to guide the use of the scale in future applied research. The motivation for this study was additionally underpinned by the growing calls for more rigorous approaches to the development and evaluation of survey measures, and iterative reconsiderations of conceptual clarity as necessary precursors for advancing educational and psychological sciences (Bringmann et al., 2022; Flake, 2021; Smaldino, 2019; van Dijk et al., 2020).

To date internal consistency reliability estimates (i.e. coefficient alpha) have been the only consistently documented evidence of scale reliability for the short-form CHAOS. Estimates of alpha are universally low (i.e.  $< .70$  in all cases) in the samples included in this study, and in other samples reported in the published literature, suggesting that the internal consistency of the items may not be sufficient for combining them in a single composite (McNeish & Wolf, 2020). In all but one of the published studies using the short-form CHAOS, the six items are either summed or averaged to compute a single composite. Creating a single composite from several items, however, assumes that a scale captures a single latent domain (McNeish & Wolf, 2020). However, previous published evidence for the short-form CHAOS (Johnson et al., 2008) and our own exploratory factor analysis suggested that the six items better represent two latent domains. An interpretational difficulty arose however, with the identification of two different patterns of item-to-factor loadings in the exploratory factor models for each of these samples. These different item-to-factor patterns, along with poor model fit in confirmatory factor models in several samples (see Table 3), support

our conclusion that the short-form CHAOS may not satisfactorily capture the quality of home environments suggested by theoretical descriptions (Bradley, 2015; Bronfenbrenner & Evans, 2000; Bronfenbrenner & Ceci, 1994; Matheny et al., 1995).

Furthermore, results of measurement invariance analyses indicated that the short-form CHAOS was non-invariant across the five included samples. Specifically, while the configural invariance model incorporating five samples was an adequate fit to the data (see Table 3, model 1.), neither the full metric invariance, nor the full scalar invariance models were acceptable according to multiple model fit criteria (Rutkowski & Svetina, 2014; West et al., 2012). The retained model was a partial metric invariance model, which constrained the factor loadings of five items to equality, and allowed the loading of one item to vary (item 5. *There is usually a television turned on somewhere in our home*). While this model was acceptable, scalar invariance was not supported, indicating that the intercepts of the items differed too greatly across samples for them to be constrained to equality. Because we included five samples in the analyses, it is neither easy nor straightforward to identify which samples and/or items were driving the scalar non-invariance – it could be a combination of any number of items across some or all samples.

The worst performing items across all samples were the two which appeared in the short-form that did not appear in the long-form version of the CHAOS: item 5. *There is usually a television turned on somewhere in our home* and item 1. *The children have a regular bedtime routine*. In terms of face validity, item 5 has become particularly dated in western cultures in the 25 years since the short-form CHAOS was proposed. For example, the data in the TEDS project was collected in the late 1990s, whereas the most recent data collection, Project KIDS, was in 2017. Compared with the 1990s, 21<sup>st</sup> century middle-class families (and children) now have access to an abundance of portable electronic devices, including smartphones, tablets, and laptops. Children and adults have access to headphones, volume control, voice-activated commands and individualized options. If the television item intended to capture ambient noise within a household, it may be outdated. If the item intended to capture parents' lack of control over children's media consumption, again, the item will likely no longer capture the range of digital media currently available to children and adolescents (Graafland, 2018).

Secondly, while the face validity of the bedtime routine item might be acceptable for samples of very young children, this routine might not be applicable to older children and adolescents. The question may capture something about personality, or developmental changes in sleep patterns (CITE), rather than an aspect of household management under the control of a parent. Unpacking the assumptions embedded in the question as it relates to variability in household order and routine raises additional questions: Is it problematic or damaging if older children lack a strictly adhered-to bedtime routine every evening? Is the amount and nature of sleep itself a better predictor of positive childhood development than regularity in bedtimes (e.g. Dewald et al., 2010)? If the CHAOS scale is to be applied in research spanning early childhood to mid-adolescence, as is currently the case, these questions, and the face validity of all the items, should be examined in the light of advancements to developmental theory, and changes to family life that have occurred since the mid-1980s when the scale was first developed in a sample of infants and toddlers.



Thirdly, the contrasting results of the two exploratory factor analyses (i.e. our own and that reported by Johnson et al., 2008) suggest that item 6 *the atmosphere in our house is calm* potentially lacks conceptual clarity (Borsboom et al., 2004; Bringmann et al., 2022). In the WRRMP data, this item loaded with others representing ‘quietness of the household’ (Johnson et al., 2008), whereas in the ADSAT data this item cross-loaded highly onto both factors, suggesting that respondents may have varying interpretations of a calm atmosphere in the home. Furthermore, whether or not a calm atmosphere is representative of poor household environments is arguable, and potentially tied to cultural norms, thus perhaps leading to the inconsistent properties of this item across datasets.

Finally, our investigation of the predictive validity of the short-form CHAOS scale was limited by the finding of measurement non-invariance. However, using two different approaches to collating variables (i.e. factor scores or averaged composites) and comparing correlations with socioeconomic status and academic achievement variables is nonetheless instructive. Using either approach to combining items, higher ratings of CHAOS correlated with poorer academic achievement in four of five samples. Factor-score correlations were generally smaller than those observed when items were averaged to create composite variables. In the WRRMP sample, CHAOS did not consistently correlate with any of the five reading or math assessments. Similarly, while SES and CHAOS were negatively correlated in general across the five samples, of note are the differences in the strength of correlations (and their statistical significance) when the average score was used rather than the factor score. In all samples, correlations using an average CHAOS composite were larger than those using the factor score. For the WRRMP and Project KIDS samples, correlations with the factor score were not significantly different from zero, whereas the correlations with the average composite were significantly different from zero. It is worth reiterating that measurement error of observed items is retained in composite variables, and can subsequently inflate or reduce covariations in unpredictable ways (Cole & Preacher, 2014; McNeish & Wolf, 2020) – as we have observed in these comparisons. While this problem can be somewhat rectified by the use of factor scores which allow differential weighting of items comprising each factor, if the observed items do not reliably capture the underlying theoretical construct, a factor score approach does not completely resolve the measurement problems (Hancock, 2003; Rhemtulla et al., 2020). The only resolution in this and many other cases is careful and considered development and renewal of items in the light of theoretical construct of interest.

### Recommendations.

Since there are multiple studies that have collected data on the short form CHAOS over the past 20 years, and several of these data sources are now accessible to researchers for secondary analyses, we provide some tentative recommendations on the use of the six CHAOS items. Given the finding that a two-factor solution is more defensible than the commonly used one-factor model we recommend that future researchers should use the two-factor configuration of the six observed items with factors termed *Noise* and *Disorder*. Results using the two-factor approach can also be compared with the one-factor approach commonly reported in the literature. We would also suggest that using factor scores, allowing items in the subscales to be differentially weighted, is more appropriate than using

sum or average scores – particularly given the variation in correlation coefficients using each approach. These recommendations, however, should not be taken as rules, and should not preclude researchers from carefully examining the properties of the items in samples not included in this study.

### Limitations

A major limitation of the analyses presented here is the non-definitive nature of the information obtained from measurement invariance tests when more than two samples are included. While we suggest that the item relating to television may be driving metric non-invariance, other items could also be contributing to this result. A second limitation is the differing ages of the children included in each sample. While we made efforts to select samples with similarly aged children, this was not always possible due to the secondary data accessed for the study. Mean age ranged from 4 years in the TEDS sample, to 11 years in both the FTP-RBE and Project KIDS samples. Differing ages of children when parents respond to the items could drive differential response patterns across samples. Nonetheless, if this is the case, it is further evidence that the short-form CHAOS is not as broadly applicable across childhood and adolescence as it is intended to be.

Finally, because the analyses are largely data-driven, the analytic choices, and the order in which different steps were undertaken in this study were affected by researcher degrees of freedom (Gelman & Loken, 2013). It would be possible to attempt different analyses and obtain different results, for example, is a one-factor solution acceptable if the television item is omitted? Or both the television and bedtime items? These different choices, however, would not get us closer to the main object of interest, which is to identify whether the six items in the short-form CHAOS are valid and reliable measure of the quality of home environments. Future work may consider these and other options within a broader program of scale development and renewal.

### Conclusion

Studies of the links between the nature of home environments and childhood development are decades old (e.g. Elardo et al., 1977; Bronfenbrenner, 1981). The original 15-item CHAOS measure clearly identified the aspects of home environments it was intended to capture. These included household disorder, high ambient noise, and lack of routine (Matheny et al., 1995) and items were developed from a wealth of earlier theorizing about how variation in home environments might relate to different aspects of early childhood development. However, the results we report here do not provide strong evidence that the short-form CHAOS adequately captures this broad and theoretically consistent construct. The rationale for selecting the six items is arguably clear: in terms of face validity and relevance, the items do cover the scope of the original construct, albeit in a more limited way. However, our findings indicate that the short form items should now be reconsidered and the scale revised in the light of more contemporary theory and contexts (e.g. Clark & Watson, 1995, 2019). In our view, a re-evaluation and update of the 15 items in the original version of the Confusion, Hubbub and Order Scale (Matheny et al., 1995) may be a useful starting point.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Funding Information

### Academic Development Study of Australian Twins

This research was supported by two Australian Research Council Discovery Project Grants: DP 120102414 (2012–2014) and DP 150102441 (2015–2018). Access to the sample was facilitated by Twins Research Australia, a national resource supported by a Centre of Research Excellence Grant (ID: 1079102), from the National Health and Medical Research Council.

### Florida Twin Project on Reading, Behaviour and Environment

This work was supported by Eunice Kennedy Shriver National Institute of Child Health & Human Development Grant 5P50HD052120–08

### Project KIDS

This work was supported by Eunice Kennedy Shriver National Institute of Child Health & Human Development Grants R21HD072286, P50HD052120, and R01HD095193.

### Western Reserve Reading and Math Project

This work was supported by Eunice Kennedy Shriver National Institute of Child Health & Human Development Grants HD38075 and HD46167

### Twins Early Development Study

TEDS is supported by a program grant from the UK Medical Research Council (MR/M021475/1 and previously G0901245), with additional support from the US National Institutes of Health (AG046938)

## Data Accessibility Statement

ADSAT data is accessible by application to the first and third authors

Available at the LDBase repository <https://ldbbase.org/> are:

- - FTB-RBE (doi:10.33009/ldbbase.1624451991.3667)
- - Project KIDS (doi: 10.33009/ldbbase.1619716971.79ee) and
- - WRRMP (doi: 10.33009/ldbbase.1643647076.d4b2)

TEDS data available by agreement with the Kings' College London <https://www.teds.ac.uk/researchers/teds-data-access-policy>

## References

- Asbury K, Dunn JF, Pike A, & Plomin R (2003). Nonshared environmental influences on individual differences in early behavioral development: A monozygotic twin differences study. *Child Development*, 74(3), 933–943. [PubMed: 12795399]
- Asbury K, Wachs TD, & Plomin R (2005). Environmental moderators of genetic influence on verbal and nonverbal abilities in early childhood. *Intelligence*, 33(6), 643–661. 10.1016/j.intell.2005.03.008
- Australian Curriculum Assessment and Reporting Authority [ACARA]. (2017). NAPLAN achievement in reading, writing, language conventions and numeracy: National report for 2017. <https://www.nap.edu.au/results-and-reports/national-reports>

- Borsboom D, Mellenbergh GJ, & van Heerden J (2004). The concept of validity. *Psychological Review*, 111(4), 1061–1071. 10.1037/0033-295X.111.4.1061 [PubMed: 15482073]
- Bradley RH (2015). Constructing and adapting causal and formative measures of family settings: The home inventory as illustration. *Journal of Family Theory & Review*, 7(4), 381–414. 10.1111/jftr.12108 [PubMed: 26997978]
- Bringmann LF, Elmer T, & Eronen MI (2022). Back to basics: The importance of conceptual clarification in psychological science. *Current Directions in Psychological Science*, 09637214221096485. 10.1177/09637214221096485
- Bronfenbrenner U (1981). *The Ecology of Human Development: Experiments by Nature and Design*. Harvard University Press.
- Bronfenbrenner U (1986). Ecology of the family as a context for human development: Research perspectives. *Developmental Psychology*, 22(6), 723–742. 10.1037/0012-1649.22.6.723
- Bronfenbrenner U, & Ceci SJ (1994). Nature-nuture reconceptualized in developmental perspective: A bioecological model. *Psychological Review*, 101(4), 568–586. 10.1037/0033-295X.101.4.568 [PubMed: 7984707]
- Bronfenbrenner U, & Evans GW (2000). Developmental science in the 21st century: Emerging questions, theoretical models, research designs and empirical findings. *Social Development*, 9(1), 115–125. 10.1111/1467-9507.00114
- Bronfenbrenner U, & Morris PA (2006). The bioecological model of human development. In Damon W, Lerner RM, & Lerner RM (Eds.), *Handbook of Child Psychology*.
- Butcher LM, & Plomin R (2008). The nature of nurture: A genomewide association scan for family chaos. *Behavior Genetics*, 38(4), 361–371. 10.1007/s10519-008-9198-z [PubMed: 18360741]
- Byrne BM (2012). *Structural Equation Modeling with Mplus: Basic Concepts, Applications and Programming*. Routledge.
- Chang L, Lu HJ, Lansford JE, Bornstein MH, Steinberg L, Chen B-B, Skinner AT, Dodge KA, Deater-Deckard K, Bacchini D, Pastorelli C, Alampay LP, Tapanya S, Sorbring E, Oburu P, Al-Hassan SM, Di Giunta L, Malone PS, Uribe Tirado LM, & Yotanyamaneewong S (2019a). External environment and internal state in relation to life-history behavioural profiles of adolescents in nine countries. *Proceedings of the Royal Society B: Biological Sciences*, 286(1917), 20192097. 10.1098/rspb.2019.2097
- Chang L, Lu HJ, Lansford JE, Skinner AT, Bornstein MH, Steinberg L, Dodge KA, Chen BB, Tian Q, Bacchini D, Deater-Deckard K, Pastorelli C, Alampay LP, Sorbring E, Al-Hassan SM, Oburu P, Malone PS, Di Giunta L, Tirado LMU, & Tapanya S (2019b). Environmental harshness and unpredictability, life history, and social and academic behavior of adolescents in nine countries. *Developmental Psychology*, 55(4), 890–903. 10.1037/dev0000655 [PubMed: 30507220]
- Cheung GW, & Rensvold RB (2002). Evaluating goodness-of-fit indexes for testing measurement invariance. *Structural Equation Modeling: A Multidisciplinary Journal*, 9(2), 233–255. 10.1207/S15328007SEM0902\_5
- Clark LA, & Watson D (1995). Constructing validity: Basic issues in objective scale development. *Psychological Assessment*, 7(3), 309–319. 10.1037/1040-3590.7.3.309
- Clark LA, & Watson D (2019). Constructing validity: New developments in creating objective measuring instruments. *Psychological Assessment*, 31(12), 1412–1427. 10.1037/pas0000626 [PubMed: 30896212]
- Coldwell J, Pike A, & Dunn J (2006). Household chaos: Links with parenting and child behaviour. *Journal of Child Psychology and Psychiatry*, 47(11), 1116–1122. 10.1111/j.1469-7610.2006.01655.x [PubMed: 17076750]
- Cole DA, & Preacher KJ (2014). Manifest variable path analysis: Potentially serious and misleading consequences due to uncorrected measurement error. *Psychological Methods*, 19(2), 300–315. 10.1037/a0033805 [PubMed: 24079927]
- Deater-Deckard K, Godwin J, Lansford JE, Tirado LMU, Yotanyamaneewong S, Alampay LP, Al-Hassan SM, Bacchini D, Bornstein MH, Chang L, Di Giunta L, Dodge KA, Oburu P, Pastorelli C, Skinner AT, Sorbring E, Steinberg L, & Tapanya S (2019). Chaos, danger, and maternal parenting in families: Links with adolescent adjustment in low- and middle-income countries. *Developmental Science*, 22(5), e12855. 10.1111/desc.12855 [PubMed: 31077512]

- Deater-Deckard K, Mullineaux PY, Beekman C, Petrill SA, Schatschneider C, & Thompson LA (2009). Conduct problems, IQ, and household chaos: A longitudinal multi-informant study. *Journal of Child Psychology and Psychiatry*, 50(10), 1301–1308. 10.1111/j.1469-7610.2009.02108.x [PubMed: 19527431]
- Dewald JF, Meijer AM, Oort FJ, Kerkhof GA, & Bögels SM (2010). The influence of sleep quality, sleep duration and sleepiness on school performance in children and adolescents: A meta-analytic review. *Sleep Medicine Reviews*, 14(3), 179–189. 10.1016/j.smrv.2009.10.004 [PubMed: 20093054]
- Dumas JE, Nissley J, Nordstrom A, Smith EP, Prinz RJ, & Levine DW (2005). Home chaos: Sociodemographic, parenting, interactional, and child correlates. *Journal of Clinical Child & Adolescent Psychology*, 34(1), 93–104. 10.1207/s15374424jccp3401\_9 [PubMed: 15677284]
- Dunn LM, & Markwardt FC (1998). Peabody Individual Achievement Test—Revised/Normative Update. Circle Pines, MN: American Guidance Service.
- Elardo R, Bradley R, & Caldwell BM (1977). A longitudinal study of the relation of infants' home environments to language development at age three. *Child Development*, 48(2), 595–603. 10.2307/1128658
- Evans GW (2006). Child development and the physical environment. *Annual Review of Psychology*, 57(1), 423–451. 10.1146/annurev.psych.57.102904.190057
- Flake JK (2021). Strengthening the foundation of educational psychology by integrating construct validation into open science reform. *Educational Psychologist*, 56(2), 132–141. 10.1080/00461520.2021.1898962
- Ganasegeran K, Selvaraj K, & Rashid A (2017). Confirmatory factor analysis of the Malay version of the Confusion, Hubbub and Order Scale (CHAOS-6) among myocardial infarction survivors in a Malaysian cardiac healthcare facility. *The Malaysian Journal of Medical Sciences: MJMS*, 24(4), 39–46. 10.21315/mjms2017.24.4.5
- Gelman A, & Loken E (2013). The garden of forking paths: Why multiple comparisons can be a problem, even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited ahead of time. Retrieved from [http://www.stat.columbia.edu/~gelman/research/unpublished/p\\_hacking.pdf](http://www.stat.columbia.edu/~gelman/research/unpublished/p_hacking.pdf)
- Gould KL, Coventry WL, Olson RK, & Byrne B (2018). Gene-environment interactions in ADHD: The roles of SES and chaos. *Journal of Abnormal Child Psychology*, 46(2), 251–263. 10.1007/s10802-017-0268-7 [PubMed: 28283857]
- Graafland JH (2018). New technologies and 21st century children: Recent trends and outcomes. OECD. 10.1787/e071a505-en
- Hancock G (2003). Fortune cookies, measurement error, and experimental design. *Journal of Modern Applied Statistical Methods*, 2(2). 10.22237/jmasm/1067644980
- Hanscombe KB, Haworth CMA, Davis OSP, Jaffee SR, & Plomin R (2010). The nature (and nurture) of children's perceptions of family chaos. *Learning and Individual Differences*, 20(5), 549–553. 10.1016/j.lindif.2010.06.005 [PubMed: 21572559]
- Hanscombe KB, Haworth CMA, Davis OSP, Jaffee SR, & Plomin R (2011). Chaotic homes and school achievement: A twin study. *Journal of Child Psychology and Psychiatry*, 52(11), 1212–1220. 10.1111/j.1469-7610.2011.02421.x [PubMed: 21675992]
- Harlaar N, Butcher LM, Meaburn E, Sham P, Craig IW, & Plomin R (2005). A behavioural genomic analysis of DNA markers associated with general cognitive ability in 7-year-olds. *Journal of Child Psychology and Psychiatry*, 46(10), 1097–1107. 10.1111/j.1469-7610.2005.01515.x [PubMed: 16178934]
- Hart SA, Petrill SA, Deater-Deckard K, & Thompson LA (2007). SES and CHAOS as environmental mediators of cognitive ability: A longitudinal genetic analysis. *Intelligence*, 35(3), 233–242. 10.1016/j.intell.2006.08.004 [PubMed: 19319205]
- Hart SA, Schatschneider C, Reynolds TR, Calvo FE, Brown BJ, Arseneault B, Hall MRK, van Dijk W, Edwards AA, Shero JA, Smart R, & Phillips JS (2020). LDBase: A Learning and Development Data Repository. 10.33009/ldbbase.

- Johnson AD, Martin A, Brooks-Gunn J, & Petrill SA (2008). Order in the house! Associations among household chaos, the home literacy environment, maternal reading ability, and children's early reading. *Merrill Palmer Quarterly*, 54(4), 445–472. 10.1353/mpq.0.0009 [PubMed: 19526070]
- Kim-Spoon J, Maciejewski D, Lee J, Deater-Deckard K, & King-Casas B (2017). Longitudinal associations among family environment, neural cognitive control, and social competence among adolescents. *Developmental Cognitive Neuroscience*, 26, 69–76. 10.1016/j.dcn.2017.04.009 [PubMed: 28544983]
- Kings College London (2022). Twins Early Development Study. <https://www.teds.ac.uk/>
- Larsen SA, Little CW, Grasby K, Byrne B, Olson RK, & Coventry WL (2020). The academic development study of Australian twins (ADSAT): Research aims and design. *Twin Research and Human Genetics*, 23(3), 165–173. 10.1017/thg.2020.49 [PubMed: 32482186]
- Lauharatanahirun N, Maciejewski D, Holmes C, Deater-Deckard K, Kim-Spoon J, & King-Casas B (2018). Neural correlates of risk processing among adolescents: Influences of parental monitoring and household chaos. *Child Development*, 89(3), 784–796. 10.1111/cdev.13036 [PubMed: 29383709]
- Laurent HK, Neiderhiser JM, Natsuaki MN, Shaw DS, Fisher PA, Reiss D, & Leve LD (2014). Stress system development from age 4.5 to 6: Family environment predictors and adjustment implications of HPA activity stability versus change. *Developmental Psychobiology*, 56(3), 340–354. 10.1002/dev.21103 [PubMed: 23400689]
- Matheny AP, Wachs TD, Ludwig JL, & Phillips K (1995). Bringing order out of chaos: Psychometric characteristics of the confusion, hubbub, and order scale. *Journal of Applied Developmental Psychology*, 16(3), 429–444. 10.1016/0193-3973(95)90028-4
- McNeish DM (2018). Thanks coefficient alpha, we'll take it from here. *Psychological Methods*, 23(3), 412–433. 10.1037/met0000144 [PubMed: 28557467]
- McNeish D, & Wolf MG (2020). Thinking twice about sum scores. *Behavior Research Methods*, 52(6), 2287–2305. 10.3758/s13428-020-01398-0 [PubMed: 32323277]
- Meredith W, & Teresi JA (2006). An essay on measurement and factorial invariance. *Medical Care*, 44(11), S69–S77. [PubMed: 17060838]
- Millsap RE, & Olivera-Aguilar M (2012). Investigating measurement invariance using confirmatory factor analysis. In Hoyle Rick H. (Ed.), *Handbook of Structural Equation Modeling* (pp. 380–392). The Guilford Press.
- Oliver BR, & Plomin R (2007). Twins' Early Development Study (TEDS): A multivariate, longitudinal genetic investigation of language, cognition and behavior problems from childhood through adolescence. *Twin Research and Human Genetics*, 10(1), 96–105. 10.1375/twin.10.1.96 [PubMed: 17539369]
- Organisation for Economic Cooperation and Development, [OECD]. (2010). TALIS 2008 Technical Report. <https://www.oecd.org/education/school/44978960.pdf>
- Petrill SA, Deater-Deckard K, Thompson LA, DeThorne LS, & Schatschneider C (2006). Reading skills in early readers: Genetic and shared environmental influences. *Journal of Learning Disabilities*, 39(1), 48–55. [PubMed: 16512082]
- Petrill SA, Pike A, Price T, & Plomin R (2004). Chaos in the home and socioeconomic status are associated with cognitive development in early childhood: Environmental mediators identified in a genetic design. *Intelligence*, 32(5), 445–460. 10.1016/j.intell.2004.06.010
- Peviani KM, Kahn RE, Maciejewski D, Bickel WK, Deater-Deckard K, King-Casas B, & Kim-Spoon J (2019). Intergenerational transmission of delay discounting: The mediating role of household chaos. *Journal of Adolescence*, 72, 83–90. 10.1016/j.adolescence.2019.03.002 [PubMed: 30875564]
- Pike A, Iervolino AC, Eley TC, Price TS, & Plomin R (2006). Environmental risk and young children's cognitive and behavioral development. *International Journal of Behavioral Development*, 30(1), 55–66. 10.1177/0165025406062124
- Plomin R, DeFries JC, Knopik VS, & Neiderhiser JM (2013). *Behavioral Genetics* (6th ed.). Worth Publisher.

- Putnick DL, & Bornstein MH (2016). Measurement invariance conventions and reporting: The state of the art and future directions for psychological research. *Developmental Review*, 41, 71–90. 10.1016/j.dr.2016.06.004 [PubMed: 27942093]
- R Core Team R (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing.
- Revelle W (2022). Psych: Procedures for psychological, psychometric, and personality research. <http://personality-project.org/r/psych/psych-manual.pdf>
- Rhemtulla M, van Bork R, & Borsboom D (2020). Worse than measurement error: Consequences of inappropriate latent variable measurement models. *Psychological Methods*, 25(1), 30–45. 10.1037/met0000220 [PubMed: 31169371]
- Rimfeld K, Malanchini M, Spargo T, Spickernell G, Selzam S, McMillan A, Dale PS, Eley TC & Plomin R (2019). Twins early development study: A genetically sensitive investigation into behavioral and cognitive development from infancy to emerging adulthood. *Twin Research and Human Genetics*, 22(6), 508–513. 10.1017/thg.2019.56 [PubMed: 31544730]
- Rosseel Y (2012). lavaan: An r package for structural equation modeling. *Journal of Statistical Software*, 48(1), 1–36. 10.18637/jss.v048.i02
- Rutkowski L, & Svetina D (2014). Assessing the hypothesis of measurement invariance in the context of large-scale international surveys. *Educational and Psychological Measurement*, 74(1), 31–57. 10.1177/0013164413498257
- Smaldino P (2019). Better methods can't make up for mediocre theory. *Nature*, 575(7781), 9. 10.1038/d41586-019-03350-5 [PubMed: 31695216]
- Smith GT, McCarthy DM, & Anderson KG (2000). On the sins of short-form development. *Psychological Assessment*, 12(1), 102–111. 10.1037//1040-3590.12.1.102 [PubMed: 10752369]
- Suku S, Soni J, Martin MA, Mirza MP, Glasgow AE, Gerges M, Van Voorhees BW, & Caskey R (2019). A multivariable analysis of childhood psychosocial behaviour and household functionality. *Child: Care, Health and Development*, 45(4), 551–558. 10.1111/cch.12665 [PubMed: 30897231]
- Taylor J, Martinez K, & Hart SA (2019). The Florida State Twin Registry. *Twin Research and Human Genetics*, 22(6), 728–730. 10.1017/thg.2019.102 [PubMed: 31685063]
- van de Schoot R, Lugtig P, & Hox J (2012). A checklist for testing measurement invariance. *European Journal of Developmental Psychology*, 9(4), 486–492. 10.1080/17405629.2012.686740
- van Dijk W, Norris CU, Otaiba SA, Schatschneider C, & Hart SA (2022). Exploring individual differences in response to reading intervention: Data from project KIDS (Kids and Individual Differences in Schools). *Journal of Open Psychology Data*, 10(1), 2. 10.5334/jopd.58 [PubMed: 36081486]
- van Dijk W, Schatschneider C, & Hart SA (2021). Open science in education sciences. *Journal of Learning Disabilities*, 54(2), 139–152. 10.1177/0022219420945267 [PubMed: 32734821]
- Wachs TD (1989). The nature of the physical microenvironment: An expanded classification system. *Merrill-Palmer Quarterly*, 35(4), 399–419.
- West SG, Taylor AB, & Wu W (2012). Model fit and model selection in structural equation modeling. In Hoyle Rick H. (Ed.), *Handbook of Structural Equation Modeling* (pp. 209–231). The Guilford Press.
- Wickham H (2014). ggplot2: An implementation of the grammar of graphics. <https://cran.microsoft.com/snapshot/2015-01-06/web/packages/ggplot2/ggplot2.pdf>
- Widaman KF, & Revelle W (2022). Thinking thrice about sum scores, and then some more about measurement and analysis. *Behavior Research Methods*. Advance online publication. 10.3758/s13428-022-01849-w
- Wilson RS, & Matheny AP (1983). Mental development: Family environment and genetic influences. *Intelligence*, 7(2), 195–215. 10.1016/0160-2896(83)90029-6
- Woodcock RW (1987). *Woodcock reading mastery tests-revised*. Circle Pines, MN: American Guidance Service

Item	Total Score Correlation
1. There is very little commotion in our home	.48
2. We can usually find things when we need them	.51
3. We almost always seem to be rushed	.55
4. We are usually able to stay on top of things	.50
5. No matter how hard we try, we always seem to be running late	.43
6. It's a real zoo in our home	.63
7. At home we can talk to each other without being interrupted	.54
8. There is often a fuss going on at our home	.60
9. No matter what our family plans, it usually doesn't seem to work out	.49
10. You can't hear yourself think in our home	.62
11. I often get drawn into other people's arguments at home	.48
12. Our home is a good place to relax	.55
13. The telephone takes up a lot of our time at home	.32
14. The atmosphere in our home is calm	.64
15. First thing in the day, we have a regular routine at home	.22

**Figure 1. Original 15-item Confusion, Hubbub and Order Scale (CHAOS) from Matheny et al. (1995)**

*Note.* Parents responded true / false to each item. Scores were summed to create a composite with higher scores representing greater 'chaos'



<p><i>Instructions:</i> Below are some things that happen in most homes. Please circle the number that best describes your home:</p> <p><i>Response options:</i> (1) Definitely untrue / (2) Somewhat untrue / (3) Not really true or untrue / (4) Somewhat true / (5) Definitely True</p>	Abbreviation
<p>Items</p> <p>1. The children / the twins / my child have (has) a regular bedtime routine* (e.g., same bedtime each night, brushing teeth, reading a story/book) (ADSAT, FTP-R, Project KIDS)</p> <p>1. The twins / the children have a regular bedtime routine* (for example, same bed each night, a bath before bed, reading a story, saying prayers) (TEDS; WRRMP)</p> <p>2. You can't hear yourself think in our home</p> <p>3. It's a real zoo in our home</p> <p>4. We are usually able to stay on top of things</p> <p>5. There is usually a television turned on somewhere in our home*</p> <p>6. The atmosphere in our house is calm</p>	<p>1. BedRoutine<sup>i</sup></p> <p>2. HomeNoise</p> <p>3. HomeZoo</p> <p>4. HomeControl<sup>i</sup></p> <p>5. HomeTV</p> <p>6. HomeCalm<sup>i</sup></p>

**Figure 2. Six items in the short-form version of the CHAOS scale with variations for different studies**

*Note.* \* indicates item that did not appear in the original 15-item scale. <sup>i</sup> indicates variables reverse-coded for analysis so that higher scores = greater 'chaos'.

**Table 1.**

Descriptive statistics of the five samples included in the analysis

Study Sample (Acronym)	Country	Wave	$N^i$	Female $^{ii}$ (%)	Age $^{iii}$		$\alpha^{iv}$	$\omega_h^v$
					M	SD		
Academic Development Study of Australian Twins (ADSAT)	Australia	1	1294	50%	8.79	0.45	0.67	0.54
Florida Twin Project on Reading, Behavior and Environment (FTP-RBE)	USA	1	568	46%	11.16	2.52	0.55	0.37
		2	437		13.30	2.44	0.63	0.53
		3	313		15.24	2.51	0.50	0.48
Project KIDS	USA	1	442	49%	11.07	3.07	0.59	0.50
Western Reserve Reading and Math Project (WRRMP)	USA	1	580	57%	6.09	0.69	0.68	0.56
		2	512		7.16	0.67	0.65	0.29
		3	494		8.21	0.82	0.70	0.63
		4	352		9.81	0.98	0.62	0.58
		5	362		10.90	1.01	0.67	0.37
		6	368		12.21	1.20	0.64	0.45
		7	246		15.05	1.45	0.59	0.48
Twins Early Development Study (TEDS)	UK	1	6009	50%	3.01	0.14	0.63	0.44
		2	8014		4.03	0.15	0.66	0.59

$^i$   $N$ =families; for twin studies the number of twins is twice the number of families.

$^{ii}$  Proportion as at study commencement.

$^{iii}$  Age calculated in years: decimal places indicate proportion of a year.

$^{iv}$  Cronbach's Alpha calculated for all six items.

$^v$  McDonald's omega (hierarchical).

Table 2

Model fit statistics testing one- and two-factor models in all samples and all waves

Sample	Model	$\chi^2$ (df)	RMSEA [90%CI]	CFI	AIC	Model Comparisons	$\chi^2$ (df)	$p$ for $\chi^2$
ADSAT	Wave 1							
	A. One-factor	126.53 (9)	0.101 [0.086, 0.117]	0.92	19821			
	B. Two-factor <i>i</i>	102.51 (8)	0.096 [0.080, 0.113]	0.94	19799	A vs B = 24.02 (1)		<.001
	C. Two-factor <i>ii</i>	<b>50.14 (8)</b>	<b>0.064 [0.048, 0.081]</b>	<b>0.97</b>	<b>19747</b>	<b>A vs C = 76.39 (1)</b>		<b>&lt;.001</b>
FTP-RBE	Wave 1							
	A. One-factor	123.12 (9)	0.150 [0.127, 0.174]	0.77	9632			
	B. Two-factor	82.04 (8)	0.128 [0.104, 0.154]	0.85	9592	A vs B = 41.08 (1)		<.001
	C. Two-factor	<b>62.59 (8)</b>	<b>0.110 [0.085, 0.136]</b>	<b>0.89</b>	<b>9573</b>	<b>A vs C = 60.53 (1)</b>		<b>&lt;.001</b>
Wave 2	A. One-factor	42.54 (9)	0.093 [0.066, 0.122]	0.92	7451			
	B. Two-factor	29.49 (8)	0.079 [0.050, 0.110]	0.95	7440	A vs B = 13.05 (1)		<.001
	C. Two-factor	<b>24.71 (8)</b>	<b>0.069 [0.039, 0.101]</b>	<b>0.96</b>	<b>7435</b>	<b>A vs C = 17.83 (1)</b>		<b>&lt;.001</b>
Wave 3	A. One-factor	34.05 (9)	0.095 [0.062, 0.130]	0.89	5327			
	B. Two-factor	29.56 (8)	0.093 [0.059, 0.131]	0.90	5325	A vs B = 4.49 (1)		.034
	C. Two-factor <i>iii</i>	<b>14.92 (8)</b>	<b>0.053 [0.000, 0.094]</b>	<b>0.97</b>	<b>5310</b>	<b>A vs C = 19.13 (1)</b>		<b>&lt;.001</b>
Project KIDS	Wave 1							
	A. One-factor	74.74 (9)	0.129 [0.103, 0.157]	0.84	7615			
	B. Two-factor	56.31 (8)	0.117 [0.090, 0.147]	0.89	7598	A vs B = 18.43 (1)		<.001
	C. Two-factor	<b>35.35 (8)</b>	<b>0.088 [0.060, 0.119]</b>	<b>0.94</b>	<b>7577</b>	<b>A vs C = 39.39 (1)</b>		<b>&lt;.001</b>
TEDS	Wave 1							
	A. One-factor	743.79 (9)	0.117 [0.110, 0.124]	0.87	101084			
	B. Two-factor	724.80(8)	0.122 [0.115, 0.130]	0.88	100657	A vs B = 18.99 (1)		<.001
	C. Two-factor	<b>313.88</b>	<b>0.080 [0.072, 0.087]</b>	<b>0.95</b>	<b>100656</b>	<b>A vs C = 429.92 (1)</b>		<b>&lt;.001</b>
Wave 2	A. One-factor	946.69 (9)	0.114 [0.108, 0.120]	0.89	133164			
	B. Two-factor	919.65 (8)	0.119 [0.113, 0.126]	0.90	133139	A vs B = 27.04 (1)		<.001
	C. Two-factor	<b>374.34 (8)</b>	<b>0.076 [0.069, 0.082]</b>	<b>0.96</b>	<b>132594</b>	<b>A vs C = 572.35 (1)</b>		<b>&lt;.001</b>
WRRMP	Wave 1							
	A. One-factor	72.57 (9)	0.110 [0.088, 0.135]	0.91	9161			
	B. Two-factor	<b>39.12 (8)</b>	<b>0.082 [0.057, 0.108]</b>	<b>0.96</b>	<b>9129</b>	<b>A vs B = 33.44 (1)</b>		<b>&lt;.001</b>
	C. Two-factor	49.47 (8)	0.095 [0.070, 0.121]	0.94	9139	A vs C = 23.09 (1)		<.001

Sample	Model	$\chi^2$ (df)	RMSEA [90%CI]	CFI	AIC	Model Comparisons	$\chi^2$ ( df)	p for $\chi^2$
Wave 2	A. One-factor	51.03 (9)	0.096 [0.071, 0.122]	0.94	7953			
	B. Two-factor	40.74 (8)	0.089 [0.063, 0.118]	0.95	7944	A vs B = 10.29 (1)		<.001
	<b>C. Two-factor</b>	<b>29.54 (8)</b>	<b>0.073 [0.046, 0.101]</b>	<b>0.97</b>	<b>7933</b>	<b>A vs C = 21.49 (1)</b>		<b>&lt;.001</b>
Wave 3	A. One-factor	113.30 (9)	0.153 [0.129, 0.179]	0.86	7718			
	B. Two-factor	103.94 (8)	0.156 [0.130, 0.183]	0.87	7711	A vs B = 9.36 (1)		.002
	<b>C. Two-factor</b>	<b>75.45 (8)</b>	<b>0.131 [0.105, 0.158]</b>	<b>0.91</b>	<b>7683</b>	<b>A vs C = 37.85 (1)</b>		<b>&lt;.001</b>
Wave 4	A. One-factor	39.69 (9)	0.098 [0.068, 0.131]	0.92	5630			
	<b>B. Two-factor</b>	<b>21.72 (8)</b>	<b>0.070 [0.035, 0.106]</b>	<b>0.96</b>	<b>5614</b>	<b>A vs B = 17.97 (1)</b>		<b>&lt;.001</b>
	C. Two-factor	31.63 (8)	0.092 [0.060, 0.126]	0.94	5624	A vs C = 8.06 (1)		.004
Wave 5	A. One-factor	52.62 (9)	0.115 [0.086, 0.146]	0.91	5754			
	<b>B. Two-factor</b>	<b>21.43 (8)</b>	<b>0.068 [0.034, 0.103]</b>	<b>0.97</b>	<b>5725</b>	<b>A vs B = 31.19 (1)</b>		<b>&lt;.001</b>
	C. Two-factor	50.64 (8)	0.121 [0.090, 0.153]	0.91	5754	A vs C = 1.98 (1)		.159
Wave 6	A. One-factor	53.66 (9)	0.116 [0.087, 0.147]	0.92	5664			
	<b>B. Two-factor</b>	<b>14.56 (8)</b>	<b>0.047 [0.000, 0.085]</b>	<b>0.99</b>	<b>5627</b>	<b>A vs B = 39.10 (1)</b>		<b>&lt;.001</b>
	C. Two-factor	48.79 (8)	0.118 [0.087, 0.150]	0.93	5661	A vs C = 4.87 (1)		.027
Wave 7	A. One-factor	16.68 (9)	0.059 [0.000, 0.102]	0.97	4036			
	<b>B. Two-factor</b>	<b>7.51 (8)</b>	<b>0.000 [0.000, 0.072]</b>	<b>1.00</b>	<b>4028</b>	<b>A vs B = 9.17 (1)</b>		<b>.002</b>
	C. Two-factor	16.43 (8)	0.065 [0.016, 0.111]	0.96	4037	A vs C = 0.26 (1)		.614

<sup>i</sup> **Model B. Two-factor** tests the model proposed by Johnson et al. (2008).

<sup>ii</sup> **Model C. Two-factor** tests the model suggested by the exploratory factor analysis of the ADSAT data.

<sup>iii</sup> This model returned negative variances (i.e. Heywood cases) for the HomeZoo and HomeCalm items.

**Table 3.**

Model fit statistics for measurement invariance tests including one wave from each of five samples

Model	$\chi^2$ (df)	RMSEA [90% CI]	CFI	AIC	Model Comparisons	$\chi^2$ ( df)	<i>p</i> for $\chi^2$
1. Configural Invariance	551.96 (40)	0.077 [0.071, 0.083]	0.957	177424			
2A. Metric Invariance	651.06 (56)	0.070 [0.065, 0.075]	0.950	177492	1 vs 2A = 99.10 (16)		<.001
<b>2B. Partial Metric Invariance</b>	<b>582.70 (52)</b>	<b>0.069 [0.064, 0.074]</b>	<b>0.955</b>	<b>177431</b>	<b>1 vs 2B = 30.74 (12)</b>		<b>.002</b>
3. Scalar Invariance <sup><i>i</i></sup>	1343.26 (68)	0.093 [0.089, 0.097]	0.892	178160	2B vs 3 = 760.57 (16)		<.001

<sup>*i*</sup>The scalar invariance model allowed for partial metric invariance – i.e. factor loadings of the TV item were allowed to vary across groups. Retained models are in bold.

**Table 4**

Factor loadings, intercepts and R-square values for each item and group for the retained partial metric invariance model

Item	ADSAT			Florida Twin Study			Project Kids			WRRMP			TEDS		
	Loading	Intercept	R <sup>2</sup>	Loading	Intercept	R <sup>2</sup>	Loading	Intercept	R <sup>2</sup>	Loading	Intercept	R <sup>2</sup>	Loading	Intercept	R <sup>2</sup>
Factor 1															
2. HomeNoise	0.86	2.63	.62	0.86	2.25	.46	0.86	1.95	.51	0.86	2.51	.70	0.86	3.27	.68
3. HomeZoo	0.92	2.13	.70	0.92	1.83	.67	0.92	1.74	.76	0.92	2.32	.73	0.92	2.66	.65
5. HomeTV <sup>1</sup>	0.35	2.72	.07	0.07	3.63	.002	0.09	3.71	.006	0.37	3.09	.09	0.47	3.26	.14
Factor 2															
1. BedRoutine	0.15	1.43	.04	0.15	1.61	.05	0.15	1.84	.03	0.15	1.46	.05	0.15	1.38	.05
4. HomeControl	0.32	1.81	.22	0.32	1.81	.19	0.32	1.94	.16	0.32	1.92	.24	0.32	1.89	.18
6. HomeCalm	0.72	2.55	.60	0.72	2.25	.76	0.72	2.11	.77	0.72	2.65	.65	0.72	2.85	.60

Note. Standardized latent factors (M=0; SD=1).

<sup>1</sup> Loadings are allowed to vary for this item, all other items loadings are constrained to equality.

**Table 5**

Correlations between one-factor and two-factor CHAOS, socioeconomic status (SES) and academic achievement criterion variables.

Study	Correlated variable	One-factor		Two-factor	
		Chaos	Noise	Disorder	Factor correlation
ADSAT	SES	-.12 ***	-.11 ***	-.09 ***	.45 ***
	Grade 3 Reading	-.11 ***	-.12 ***	-.04	
	Grade 3 Math	-.07 *	-.07 *	-.03	
FTP-RBE Wave 1	SES	-.10 *	-.15 ***	-.09 *	.21 ***
	FCAT Reading 2011–12	-.16 *	-.26 **	-.06	
	FCAT Reading 2012–13	-.21 ***	-.27 ***	-.08	
Project KIDS	SES	-.09	-.23 ***	-.07	.29 ***
	English Language Arts Grade	-.18 ***	-.20 ***	-.20 ***	
	Math Grade	-.06	-.07	-.17 ***	
WRRMP Wave 4 <sup>i</sup>	SES <sup>ii</sup>	-.04	-.03	.01	.20 ***
	PIAT passage comprehension	.01	.00	.01	
	WRMT passage comprehension	-.04	-.05	-.00	
	WJ Calculation	-.06	-.06	-.03	
	WJ applied problems	-.01	-.02	-.04	
	WJ quantitative concepts	-.03	-.04	.00	
TEDS Wave 4	SES	-.23 ***	-.24 ***	-.05 ***	.33 ***
	English National Curriculum Assessment <sup>iii</sup>	-.16 ***	-.16 ***	-.07 *	
	Math National Curriculum Assessment <sup>iii</sup>	-.11 ***	-.12 ***	-.04 *	

Note. Correlations are between factors and criterion variables.

\*  $p < .05$ ;

\*\*  $p < .01$ ;

\*\*\*  $p < .001$ .

<sup>i</sup>For the WRRMP Wave 4 data we use the same configuration of items as the remaining datasets for the two factor models, notwithstanding the better fit of the alternative model. Interestingly correlations remained non-significant with the alternative item configuration reported in Johnson et al. (2008).

<sup>ii</sup>SES proxy variable comprising an average of both parents' educational attainment.

<sup>iii</sup>English and math assessments at age 7.

**Table 6**

Correlations between CHAOS average composites, socioeconomic status (SES) and academic achievement criterion variables

Study	Correlated variable	One variable		Two variables	
		Chaos	Noise	Disorder	Composite score correlation
ADSAT	SES	-.21***	-.23***	-.10***	.41***
	Grade 3 Reading	-.16***	-.20***	-.04	
	Grade 3 Math	-.12***	-.15***	-.04	
FTP-RBE Wave 1	SES	-.21***	-.23***	-.10*	.27***
	FCAT Reading 2011–12	-.24***	-.30***	-.06	
	FCAT Reading 2012–13	-.25***	-.29***	-.10	
Project KIDS	SES	-.21***	-.25***	-.07	.31***
	English Language Arts Grade	-.26***	-.22***	-.20***	
	Math Grade	-.15***	-.07	-.18***	
WRRMP Wave 4 <sup>i</sup>	SES <sup>ii</sup>	-.15*	-.19***	-.03	.39***
	PIAT passage comprehension	-.09	-.12*	-.01	
	WRMT passage comprehension	-.16***	-.20***	-.04	
	WJ Calculation	-.10	-.13*	-.03	
	WJ applied problems	-.10	.14*	.00	
	WJ quantitative concepts	-.07	-.12	.02	
TEDS Wave 4	SES	-.32***	-.38***	-.09***	.36***
	English National Curriculum assessment <sup>iii</sup>	-.21***	-.23***	-.09***	
	Math National Curriculum Assessment <sup>iii</sup>	-.15***	-.17***	-.06***	

Note. Composite variables created by averaging six items for one-factor CHAOS; three items each for *Noise* and *Disorder*.

<sup>ii</sup> SES proxy variable in this dataset is an average of the educational attainment of both parents.

<sup>iii</sup> English and math assessments at age 7.