

# Development of Internalizing Problems From Adolescence to Emerging Adulthood: Accounting for Heterotypic Continuity With Vertical Scaling

Isaac T. Petersen  
University of Iowa

Oliver Lindhiem  
University of Pittsburgh

Brandon LeBeau  
University of Iowa

John E. Bates  
Indiana University

Gregory S. Pettit  
Auburn University

Jennifer E. Lansford and Kenneth A. Dodge  
Duke University

Manifestations of internalizing problems, such as specific symptoms of anxiety and depression, can change across development, even if individuals show strong continuity in rank-order levels of internalizing problems. This illustrates the concept of heterotypic continuity, and raises the question of whether common measures might be construct-valid for one age but not another. This study examines mean-level changes in internalizing problems across a long span of development at the same time as accounting for heterotypic continuity by using age-appropriate, changing measures. Internalizing problems from age 14–24 were studied longitudinally in a community sample ( $N = 585$ ), using Achenbach's Youth Self-Report (YSR) and Young Adult Self-Report (YASR). Heterotypic continuity was evaluated with an item response theory (IRT) approach to vertical scaling, linking different measures over time to be on the same scale, as well as with a Thurstone scaling approach. With vertical scaling, internalizing problems peaked in mid-to-late adolescence and showed a group-level decrease from adolescence to early adulthood, a change that would not have been seen with the approach of using only age-common items. Individuals' trajectories were sometimes different than would have been seen with the common-items approach. Findings support the importance of considering heterotypic continuity when examining development and vertical scaling to account for heterotypic continuity with changing measures.

*Keywords:* internalizing problems, heterotypic continuity, vertical scaling, changing measures, developmental trajectories

*Supplemental materials:* <http://dx.doi.org/10.1037/dev0000449.supp>

Internalizing problems, including depression and anxiety, are among the most common and burdensome problems that adolescents and adults experience. The broadband, dimensional concept of internalizing problems, which represents multiple, specific symptoms, was discovered through factor analytic test development work (Achenbach & Edelbrock, 1978). Key findings with

this concept include the following: (a) the internalizing problems dimension captures individuals' meaningful clinical and subclinical difficulties, including risk for anxiety or mood disorders, relationship conflicts, ineffective parenting, and poor health (e.g., Eaton et al., 2013); (b) internalizing problems have well-established norms; and (c) they show both considerable stability in individual

---

*Editor's Note.* Elliot M. Tucker-Drob served as the action editor for this manuscript—EFD

---

This article was published Online First November 20, 2017.  
Isaac T. Petersen, Department of Psychological and Brain Sciences, University of Iowa; Oliver Lindhiem, Department of Psychiatry, University of Pittsburgh; Brandon LeBeau, Department of Psychological and Quantitative Foundations, University of Iowa; John E. Bates, Department of Psychological and Brain Sciences, Indiana University; Gregory S. Pettit, Department of Human Development and Family Studies, Auburn University; Jennifer E. Lansford and Kenneth A. Dodge, Center for Child and Family Policy, Duke University.

---

The authors acknowledge Olga Berkout and Rachel Salk for their helpful feedback on a draft of this article. The Child Development Project has been funded by Grants MH42498, MH56961, MH57024, and MH57095 from the National Institute of Mental Health; HD30572 from the Eunice Kennedy Shriver National Institute of Child Health and Human Development; and DA016903 from the National Institute on Drug Abuse. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Institutes of Health.

Correspondence concerning this article should be addressed to Isaac T. Petersen, Department of Psychological and Brain Sciences, University of Iowa, 309 Stuit Hall, Iowa City, IA 52242. E-mail: [isaac-t-petersen@uiowa.edu](mailto:isaac-t-petersen@uiowa.edu)

differences over time and some rank-order change along with (d) some normative (mean-level) fluctuations across development. These findings suggest further questions toward understanding how internalizing problems develop and toward charting both normative patterns and individual trajectories of internalizing problems.

According to considerable prior research, rates of depression and other internalizing problems have been shown to increase in adolescence, peak in mid-to-late adolescence, and decrease into adulthood (Adkins, Wang, & Elder, 2009). As an example of robust sex differences in the development of internalizing problems, females show higher levels of depression than males, with sex differences emerging around the onset of puberty and the greatest sex differences appearing in mid-to-late adolescence (Hankin et al., 1998). The marked developmental changes and sex differences in depression during this developmental era make the transition from adolescence to adulthood particularly important to study.

### Heterotypic Continuity

Among the many studies of internalizing problems, we have found none that examined trajectories of internalizing problems using measures adjusted to maintain construct validity and consider mean-level change over a lengthy developmental span. This is important because it appears that internalizing problems change in their manifestation over time (Avenevoli & Steinberg, 2002) and different measures may be needed at different ages to accurately understand how internalizing problems develop. Internalizing problems may manifest differently in adolescents compared to adults. It has been shown that somatic complaints (e.g., headaches, stomachaches, heart pounding) are more strongly associated with and more common in those with internalizing problems earlier than later in development (Achenbach, 1991, 1997; Ryan et al., 1987). This is an example of *heterotypic continuity*, which refers to persistence of an underlying construct or process with manifestations that change over the course of development (Petersen, Hoyniak, McQuillan, Bates, & Staples, 2016). Heterotypic continuity occurs when the *same* psychological reasons underlie *different* behaviors at different ages. Heterotypic continuity is analogous to the transformation of water to ice or steam, or of a caterpillar to a butterfly—the underlying core is preserved but the manifestation changes. Using measures that change with development to maintain construct validity of internalizing problems could be important for better understanding of (a) the normative trajectory of internalizing problems across ages, (b) individual differences in trajectories of internalizing problems, and (c) how risk and protective factors influence individuals' development of internalizing problems. This would build construct validity and advance understanding of how internalizing problems develop. Developmental psychology seeks to understand processes of continuity and change across the life span and not just limited windows of time or stages of life. However, studying heterotypic continuity over long spans of development poses methodological and theoretical challenges and opportunities.

### The Challenge of Heterotypic Continuity When Examining Development

**Measuring the development of internalizing problems.** Because of the heterotypic continuity of internalizing problems, measures have been designed to accommodate changes in the manifestation of internalizing problems. Most notably, the inter-

nalizing scale of the Youth Self-Report (YSR; Achenbach, 1991) designed for 11- to 18-year-olds includes items reflecting anxiety, depression, and somatic complaints. The internalizing scale on the Young Adult Self-Report (YASR; Achenbach, 1997) designed for 18- to 30-year-olds includes items reflecting anxiety and depression but not somatic complaints. To chart internalizing problems from adolescence to adulthood using the YSR and YASR, then, it is a challenge to measure participants' actual change in internalizing problems, despite changing measures.

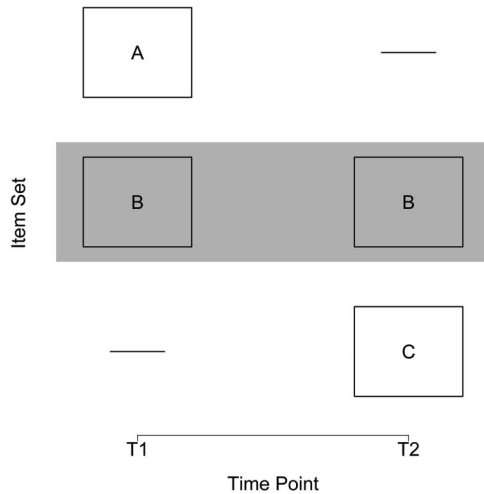
**Heterotypic continuity as the focus of study.** Heterotypic continuity is a developmental phenomenon examined by many researchers using structural equation modeling (SEM) or item response theory (IRT). Examining how strongly items relate to a latent trait (i.e., item factor loadings in SEM or item discrimination in IRT) can help determine which behaviors most strongly reflect a construct at a given point in development (i.e., continuity of the factor structure at the behavior/item level). Other researchers have examined the continuity of constructs at the latent/syndrome level (e.g., Snyder, Young, & Hankin, 2017). Given how much developmental research demonstrates how constructs change in their expression over time, surprisingly little research in developmental psychology has explored the best ways to account for heterotypic continuity, that is, how to examine individuals' developmental trajectories in constructs that change in their manifestation over time.

**Accounting for the heterotypic continuity of internalizing problems.** When examining continuity and change of individuals' trajectories on a construct, heterotypic continuity can be useful in advancing developmental theory and practice. It can also become a confound that needs to be accounted for, rather than the focus of study. When examining change, especially over a lengthy developmental span, it is important to consider and, if necessary, account for heterotypic continuity. If heterotypic continuity is not properly accounted for, the same measure may not reflect the same construct across time and, therefore, scores on the measure may not be comparable across time. To account for heterotypic continuity, changes in measurement should accommodate changes in the manifestation of the construct to retain construct validity invariance (Knight & Zerr, 2010). For example, for developmental reasons, the measurement of internalizing problems should assess somatic problems to a greater degree earlier in development. Thus, the consequence of heterotypic continuity is that *different* items over time may be necessary to assess the *same* construct over time. There are three primary approaches to measuring a construct over time, each with its own advantages and limitations.

### Approaches to Measuring a Construct Over Time

The three approaches to measuring a construct over time include administering (a) all possible items across all ages, (b) only the common items across all ages, and (c) the construct-valid items at each age. Traditionally, developmental psychologists have used all possible items (Approach 1) or only the common items (Approach 2) across all ages when measuring a construct over time. Below, we discuss the approaches and the importance of using the construct-valid items at each age instead (Approach 3), which is depicted in Figure 1.

**All possible items across all ages.** The first approach to measuring a construct over time uses all possible items across all



*Figure 1.* Depiction of using construct-valid items at each age with a common-item design. Item set A corresponds to items that are construct-valid at only T1. Item set B corresponds to items that are construct-valid at both T1 and T2. Item set C corresponds to items that are construct-valid at only T2. The “common items” (item set B) are highlighted in gray. The present study used the construct-valid items at each age (Approach 3: i.e., item sets A and B at T1 and item sets B and C at T2), by using the common items to link the different item sets. Although there were more than two time points (10) in the present study, we used only two different measures (hence we depict the common item-design with T1 and T2 for simplicity).

ages. One advantage is that it is a comprehensive approach to assessment that allows examining change in each item across the developmental span. The approach has key disadvantages, however. First, it is inefficient. It requires extra time to assess all items across all ages. Second, it could assess items that are developmentally inappropriate at a given age (because of changes in item difficulty or severity, i.e., how infrequently an item is correct or endorsed). For example, in a test of math ability, it would be developmentally inappropriate to ask a 7-year-old an advanced calculus question. Third, the aggregation of scores on all possible items could result in a score that lacks construct validity invariance and therefore becomes incomparable over time if the construct changes in its manifestation (because of changes in item discrimination, i.e., how strongly the item relates to the trait). For example, a measure including somatic complaints to assess internalizing problems in adulthood may not reflect the same construct as assessed by the same measure of internalizing problems in adolescence. Thus, the same measure may not reflect the same construct at different ages. Thus, aggregating scores on all possible items across all ages could produce problems for interpretation when heterotypic continuity is likely.

**Only the common items across all ages.** The second approach to measuring a construct over time is to use only the common items across all ages. Using only the common items across all ages has the advantage that it is efficient, but also has key disadvantages. First, using only the common items results in a loss of information because there are fewer items assessing the construct, which may make the measure less sensitive to developmental change. Excluding from all ages items that are developmentally inappropriate at some ages and developmentally appro-

priate at other ages could result in the systematic loss of information on the full scope of internalizing problems, which is crucial for assessing individual differences, especially at low and clinical levels of problems. For instance, in a hypothetical study of internalizing problems from 2 to 18 years of age, suicidality would not likely be used as a common item across all ages because it would be developmentally inappropriate to ask parents of a 2-year-old whether the child is suicidal. Omitting suicidality across all ages, however, would result in a loss of information regarding an important internalizing problem with high severity. Second, the measure may lack content validity because it is not measuring the construct as a whole, in particular, the age-specific manifestations.

**Construct-valid items at each age.** The third approach to measuring a construct over time is to use the construct-valid items at each age. In the context of heterotypic continuity, this would mean using different items across time—those items at a given age that are valid for the target construct. Using the construct-valid items at each age has several advantages. First, it retains content validity and construct validity invariance. Second, it is more efficient than using all possible items across all ages. And where there is heterotypic continuity, it is the best way to maintain construct validity invariance.

There are still important issues in using the construct-valid items at each age when the items differ across time. For one, there is the question of how to measure individuals’ change in a construct when a different measure is used at each age. Different scores on the measures’ different items over time could reflect either (a) a person’s change in the trait, or (b) an artifactual change resulting from the different measures/items at each age having different meaning. Assuming the measures reflect the same construct over time (i.e., construct validity invariance), the next consideration for determining whether different scores for an individual over time reflect actual change is the issue of statistical equivalence. Are the measures’ scores on the same metric or scale so they can be meaningfully compared? First, the measures should have the same range of possible scores. Second, in order to measure absolute change (rather than solely an individual’s change relative to others), a score on the measure at Time 1 (T1) should reflect the same trait level on the construct as the same score on the measure at T2. There are several possible solutions to ensuring statistical equivalence of different measures over time, including: (a) age-norming, (b) average/percentage scores, and (c) vertical scaling.

**Age-norming.** Age-norming (e.g., standard scores and percentiles) is commonly used to compare scores on different measures because age-normed scores have a similar mathematical metric. Standard scores (e.g., *t*- or *z*-scores) have a fixed mean and standard deviation. Percentiles have a fixed range (0–100). Age-norming can be useful for examining individuals’ *relative* change (i.e., change relative to other individuals in the sample or relative to a norm-referenced sample). However, because age-normed scores have a fixed scale, they cannot detect *absolute* change (i.e., change in an individual’s trait level or the group mean or variability over time). Standardizing scores with a fixed range or mean and standard deviation does not ensure the scores are on the same metric, so age-norming is generally inadvisable when examining development (Moeller, 2015).

**Average or percentage scores.** Another approach to comparing scores across different measures is an average score or percentage score that accounts for the different number of items in

each measure. A major assumption of average and percentage scores is that the items on the different measures do not differ in discrimination or severity (defined in the next paragraph). However, it is unlikely that measures with different items will have the same severity, especially when the item content differs across the two measures. Thus, average or percentage scores are not advisable in most contexts dealing with different measures over time. To compare scores on different measures over time, researchers recommend vertical scaling (Kolen & Brennan, 2014).

**Vertical scaling.** In vertical scaling, measures that assess a similar construct but differ in difficulty or severity are placed on the same scale. Vertical scaling is widely used in educational testing because the same test items tend to become easier relative to a given level of ability as children get older. Multiple approaches exist for vertical scaling. For the present study, we used the IRT approach to vertical scaling (Kolen & Brennan, 2014). We fit IRT models that estimate two properties of each item: (a) discrimination and (b) difficulty (severity). An item's discrimination parameter describes how well the item distinguishes between low and high levels of the trait. For example, an item asking how often a person feels depressed will have a higher discrimination for internalizing problems compared to an item asking how often the person reads. An item's difficulty parameter describes the trait level (level on the latent criterion) at which the probability of endorsing the item is 50%. In the context of psychopathology, a higher difficulty parameter reflects a higher, more severe level of internalizing problems, so henceforth we refer to the difficulty parameter as severity. For example, an item asking how often a person thinks about suicide will have a higher severity compared to an item asking how often the person feels sad. Based on the items' parameters and participants' responses on the items, IRT estimates each person's latent trait level of internalizing problems (i.e., ability score or theta).

When the different measures have common items over time, the IRT approach to vertical scaling uses common items administered across ages to link the measures on the same scale by finding scaling parameters that put the trait level scores on the same metric. The scaling parameters are determined as the linear transformation (i.e., the intercept and slope parameter) that, when applied to the second measure, minimizes the differences between the probability of an individual endorsing the common items across the two measures. Although the common items are used to determine the general form of change on the same scale, all developmentally relevant, construct-valid items are used to estimate each person's trait level on this scale.

A number of studies have used vertical scaling in the fields of education and cognitive testing to measure growth with changing measures over time. As one prime example, McArdle et al. (2009) examined the development of cognitive ability from 2 to 72 years of age.

### Limitations of Previous Research

Despite the numerous studies using vertical scaling in education and cognitive testing, to our knowledge, no studies have examined the development of psychopathology or social development more generally using vertical scaling. Moreover, despite researchers acknowledging the importance of examining the heterotypic continuity of internalizing problems (Sterba, Prinstein, & Cox, 2007),

to our knowledge, no studies have examined trajectories of internalizing problems with changing measures to account for heterotypic continuity, maintain construct validity, and examine mean-level change over a lengthy span of development.

To our knowledge, the only studies examining trajectories of *broadband* internalizing problems with changing measures come from the Australian Temperament Project, which examined trajectory classes from ages 3–15 (and anxiety and depression from ages 11–27; Betts et al., 2016; Letcher, Sanson, Smart, & Toumbourou, 2012; Letcher, Smart, Sanson, & Toumbourou, 2009; Toumbourou, Williams, Letcher, Sanson, & Smart, 2011). The studies did not link the different measures or account for changes in the measures' scales, however, so they did not allow interpretation of mean-level change across measurement changes.

The challenge of heterotypic continuity has led researchers to frequently grapple with the issue of developmental equivalence or to avoid examining the development of internalizing problems across lengthy spans. Many studies have used all possible items or only the common items to maintain the same measure over time. Regarding all possible items, we have seen many studies of internalizing problems that have used measures outside the ages they were originally designed to assess (Adkins et al., 2009; Broeren, Muris, Diamantopoulou, & Baker, 2013; Côté et al., 2009; Crocetti, Klimstra, Keijsers, Hale, & Meeus, 2009; Hale, Raaijmakers, Muris, van Hoof, & Meeus, 2008; Leadbeater, Thompson, & Gruppuso, 2012; Mathiesen, Sanson, Stoolmiller, & Karevold, 2009; Meadows, Brown, & Elder, 2006; Miers, Blöte, de Rooij, Bokhorst, & Westenberg, 2013; Morin et al., 2011). We have also seen studies of internalizing problems that used only common items over time (Fanti & Henrich, 2010; Gilliom & Shaw, 2004; Sterba et al., 2007).

Thus, studies frequently deal with the issue of developmental equivalence and, in many cases, resort to using a measure at an age outside the age range of validation or to discarding relevant items. An important advance for the field is learning how to handle changes in measurement, because ignoring the heterotypic continuity of internalizing problems over lengthy spans of development likely results in measures that violate construct validity (if summing all possible items) or content validity (if using only common items). Moreover, it also allows measuring internalizing problems in age-appropriate ways, for the sake of understanding development across important developmental transitions.

To ignore heterotypic continuity by using only common items and discarding items reflecting the age-specific manifestations of internalizing problems (e.g., somatic complaints) results in a measure that captures the development of specific problems (i.e., the common items) without capturing the development of the *construct* of internalizing, and can result in inaccurate trajectories. Chen and Jaffee (2015) found that the common items failed to detect the adolescent-onset of externalizing problems observed in a subgroup when using age-relevant items. The challenge of heterotypic continuity may account for why we have not seen studies examining trajectories of broadband internalizing problems across the transition from adolescence into adulthood. We approached this problem by comparing different approaches to vertical scaling. Vertical scaling approaches are widely used in other fields to examine change with different measures over time (Kolen & Brennan, 2014), so it seemed plausible that vertical scaling would



be a useful approach to account for heterotypic continuity in developmental psychology.

### The Present Study

We examined the development of internalizing problems over a decade of life, and used vertical scaling with different measures over time because internalizing problems demonstrate heterotypic continuity. After rescaling the different measures of internalizing problems to be on the same scale to account for heterotypic continuity, we examined growth curves of internalizing problems and whether the trajectories differed by sex or ethnicity.

### Method

#### Participants

Children ( $N = 585$ ) were recruited for the Child Development Project (Dodge, Bates, & Pettit, 1990) from two cohorts in 1987 and 1988 from schools at three sites: Nashville, TN; Knoxville, TN; and Bloomington, IN. The schools and the sample represented families with a broad range of socioeconomic status (SES), representative of the populations at the respective sites. The Hollingshead index of SES ( $M = 39.53$ ,  $SD = 14.01$ , range: 8 to 66, stratum 1: 17% of the sample, 2: 33%, 3: 25%, 4: 16%, 5: 9%) reflected a broad range for the original sample, which was 52% male, 81% European American, 17% African American, and 2% of “other” ethnicity. Over the course of the project, the Child Development Project protocols have been approved by Institutional Review Boards (IRBs) at Indiana University, Vanderbilt University, the University of Tennessee, Auburn University, and Duke University. The current protocol “How Chronic Conduct Problems Develop” (protocol number 40) is approved by the Duke University IRB.

Children were followed annually with parents’, teachers’, peers’, and/or self-report ratings of the children’s internalizing problems. The present study focuses on self-report ratings of adolescents’ and young adults’ internalizing problems from 14 to 24 years of age. We focused on self-reports because of the accuracy of adolescents’ reports of their own internalizing problems—adolescents are in a unique position to report on their subjective experiences of internalizing problems (De Los Reyes & Kazdin, 2005).

#### Measures

Adolescents rated their level of internalizing problems annually on the Internalizing Scale of the YSR (Achenbach, 1991) from ages 14 to 19 (except age 18). From ages 20 to 24, they rated their internalizing problems annually on the Internalizing Scale of the YASR (Achenbach, 1997). Adolescents rated internalizing problems on the YSR and YASR as *not true*, *somewhat or sometimes true*, or *very true or often true*, scored 0, 1, and 2, respectively (although vertical scaling does not require scores on the same response scale). Scores on the internalizing scale were summed across items. Internal consistency of items ranged from  $\alpha = .89$  to  $.91$ , depending on the year. The Achenbach scales have strong validity, including content, construct, and criterion-related validity (Sattler & Hoge, 2006).

Items on the internalizing scale differed somewhat between the YSR (31 items) and YASR (23 items) in ways that reflected the heterotypic continuity of internalizing problems. For instance, somatic complaints were included in the measure of adolescents’ internalizing problems on the YSR, but they were not included in the measure of young adults’ internalizing problems on the YASR. The YSR internalizing scale included withdrawn, somatic complaints, and anxious/depressed subscales, whereas the YASR internalizing scale included withdrawn and anxious/depressed subscales. The internalizing scale on the YSR and YASR shared 17 common items, while 14 of the items on the Internalizing scale were unique to the YSR and six were unique to the YASR.<sup>1</sup> Descriptive statistics and a Pearson correlation matrix of the raw internalizing sum scores at each age are in Table 1. Possible scores ranged from 0–62 on the YSR, 0–46 on the YASR, and 0–34 on the 17 common items of the YASR and YSR, with higher scores reflecting higher levels of internalizing problems.

We also examined the association of scores on the internalizing scale with scores on the internalizing and externalizing scales 1 year later. Internalizing problems showed strong convergent and discriminant validity. The average correlation of internalizing problems with later internalizing problems was  $r = .72$  (95% CI [.67, .76]). The average correlation of internalizing problems with later externalizing problems was  $r = .41$  (95% CI [.33, .49]).

#### Statistical Analysis

In the present study, we used the IRT approach to vertical scaling (as described in Kolen & Brennan, 2014) to transform scores on the YASR to the scale of the YSR. In the context of vertical scaling, IRT estimates people’s latent trait scores of internalizing problems over time (i.e., a latent variable or “true score” approach to vertical scaling). As further validation of the findings from the IRT approach to vertical scaling, we also conducted an alternative approach to vertical scaling, known as Thurstone scaling (as described in Kolen & Brennan, 2014). Unlike IRT, the Thurstone scaling approach to vertical scaling retains the raw metric (i.e., an observed score or “raw score” approach to vertical scaling). Our findings from the Thurstone scaling approach are available in supplementary Appendix S3.

**IRT models.** Internalizing problems were analyzed with graded response models in IRT using the mirt package (Chalmers, 2012) in R. The mirt package uses an expectation-maximization algorithm known as marginal maximum likelihood, which uses all available data and provides valid inferences when data are missing at random or completely at random. Graded response models allow polytomous variables with more than two response categories (e.g., 0–2 Likert scale in the present study). The models estimated three parameters for each item: (a) discrimination ( $a$ ); (b) severity for the threshold from 0–1 ( $b_1$ ); and (c) severity for the threshold from 1–2 ( $b_2$ ). We examined model fit with RMSEA and CFI. We fit a separate IRT model at each age for the purposes of linking the measures across time, rather than fitting all items in the same model (i.e., concurrent calibration). Although concurrent calibration procedures tend to have greater precision of item parameter estimates, separate estimation is

<sup>1</sup> The YASR item reflecting whether the adult was concerned about his or her looks was not administered at each age the YASR was administered, so it was not included in our calculations of the internalizing scale.

Table 1  
Pearson Correlation Matrix (Two-Tailed) of Raw Internalizing Problem Scores and Descriptive Statistics

Age	14	15	16	17	19	20	21	22	23	24
14	—									
15	.68	—								
16	.59	.68	—							
17	.58	.59	.65	—						
19	.50	.53	.62	.68	—					
20	.39	.45	.53	.62	.68	—				
21	.38	.42	.46	.61	.63	.72	—			
22	.37	.38	.44	.58	.60	.70	.75	—		
23	.39	.45	.50	.59	.62	.69	.72	.82	—	
24	.41	.41	.45	.54	.60	.67	.67	.76	.80	—
<i>n</i>	412	407	452	429	464	479	465	466	486	464
Missing %	30	30	23	27	21	18	21	20	17	21
<i>M</i>	9.44	9.90	9.63	9.03	8.66	8.52	8.71	8.94	8.96	9.45
<i>SD</i>	7.42	7.92	7.52	7.73	7.03	6.83	6.97	6.93	7.07	7.50

Note. All correlations are significant at  $p < .001$  level. Dashed lines separate the scores from the Youth Self-Report (YSR; ages 14–19) from the Young Adult Self-Report (YASR; ages 20–24). Note that the mean scores on the YSR versus YASR are not directly comparable on the same metric because they had different numbers (and types) of items in the calculation of the Internalizing scale (YSR: 31 items, YASR: 23 items).

considered safer over lengthy developmental spans because the unidimensionality assumption of IRT is more likely to be violated in concurrent calibration (Kolen & Brennan, 2014).

**Vertical scaling.** Vertical scaling involves placing two measures that assess a similar construct but differ in difficulty/severity on the same scale. Ideally, the two measures should have some items with the same contents to ensure scores on the measures can be linked (i.e., made comparable). In the present study, we used the IRT approach to vertical scaling to transform scores on the YASR to the scale of the YSR. The YSR and YASR have different but overlapping item content, so we needed to put them on the same scale. We applied vertical scaling that scales the scores across the different measures using the items that are in common across both measures (i.e.,

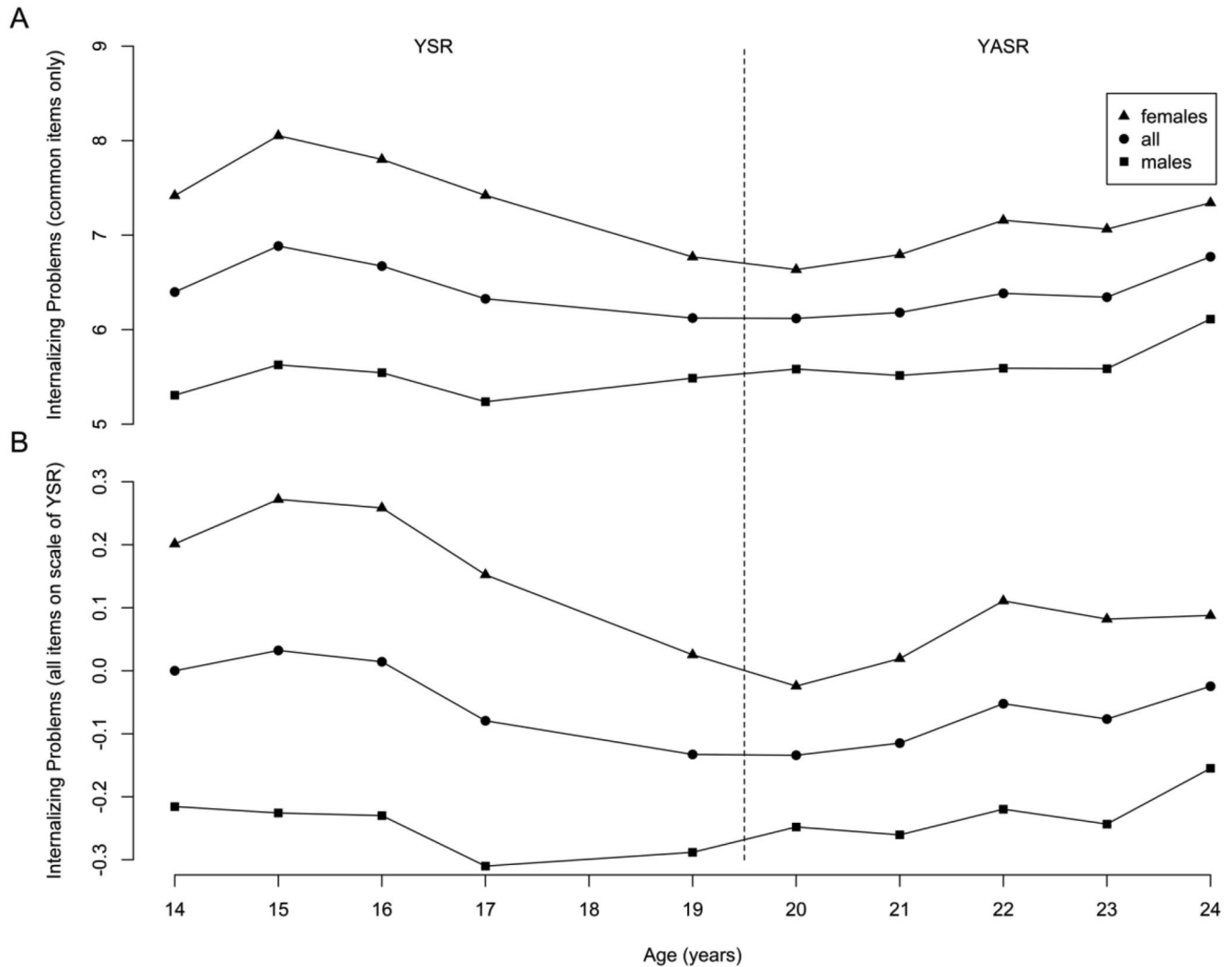
a common-item nonequivalent group or anchor instrument design, see Figure 1). We applied the following steps for vertical scaling in a common-item design (Kolen & Brennan, 2014) to link YASR scores with YSR scores:

1. To ensure a meaningful mean-level of change of internalizing problem scores across ages 14–24, we first examined scores on the 17 common items (i.e., the items that were common to both the YSR and YASR). Participants’ mean scores on the common items are in Table 2 and are depicted in Panel A of Figure 2.
2. As described earlier, we fit separate IRT models at each age.

Table 2  
Descriptive Statistics and Scaling Parameters of the Vertically Scaled Scores, Along With Descriptive Statistics of the Common Items

Age	Common items		Vertically scaled scores		Scaling parameters	
	<i>M</i>	<i>SD</i>	<i>M</i>	<i>SD</i>	A	B
14	6.398	4.897	.000	1.000	—	—
15	6.885	5.437	.032	1.066	1.091	.017
16	6.673	5.209	.014	1.012	.925	-.003
17	6.326	5.284	-.079	1.046	1.064	-.113
19	6.123	4.954	-.133	1.004	.950	-.043
20	6.119	5.164	-.134	1.035	1.048	-.012
21	6.181	5.243	-.115	1.045	1.001	.021
22	6.384	5.252	-.052	1.019	.960	.070
23	6.344	5.313	-.077	1.064	1.068	-.039
24	6.772	5.606	-.025	1.115	1.057	.041

Note. The vertically scaled scores are rescaled to be on the reference scale of the YSR at age 14. The scaling parameters are calculated in reference to the previous year. For example, the age 16 scaling parameters reflect the scaling parameters to link age 16 to age 15. To link age 16 to the reference scale at age 14, however, a process of linking and chaining is necessary (linking age 16 to age 15 using the age 16 scaling parameters, and then chaining them to age 14 by linking age 15 to age 14 using the age 15 scaling parameters). For instance, trait level scores at age 15 were rescaled to the scale at age 14 by multiplying the age 15 trait level scores by 1.091 and adding .017. Trait level scores at age 16 were rescaled to the scale at age 14 by first multiplying the age 16 trait level scores by .925 and subtracting .003 to put them on the age 15 scale, and then multiplying the new scores by 1.091 and adding .017 to put them on the age 14 scale.



*Figure 2.* Panel A depicts participants' mean raw scores on the common items (i.e., the items that were common to the internalizing scale of the Youth Self-Report [YSR] and Young Adult Self-Report [YASR]). Panel B depicts participants' mean internalizing problem scores on *all* age-relevant items of the internalizing scale, after rescaling the YASR (and YSR) scores to the metric of the YSR (based on the IRT metric of the YSR at age 14). Internalizing problems to the left of the dashed line (i.e., ages 14–19) were rated on the YSR. Internalizing problems to the right of the dashed line (i.e., ages 20–24) were rated on the YASR. Internalizing problem reports were not collected at age 18.

3. We used vertical scaling procedures to calculate scaling parameters that linked the IRT factor scores (trait level scores or theta) from the YSR and YASR at different ages on the same scale. Vertical scaling uses common items administered across ages to link the measures on the same scale by finding scaling parameters that put the trait level scores on the same metric. We used the *plink* package (Weeks, 2010) in R to calculate Stocking-Lord scaling parameters. To link the measures, scaling parameters were calculated using an iterative algorithm that minimizes the sum of squared differences between the expected aggregate scores for the common items for each measure. Thus, the scaling parameters minimize the dif-

ferences between the probability of an individual endorsing the common items across the two measures (or ages).

To calculate scaling parameters, we first set the target scale to be the YSR at age 14, and calculated scaling parameters at age 15 to link the YSR scores at age 15 to be on the same scale as the YSR at age 14. We then applied a process of linking and chaining (Kolen & Brennan, 2014) to calculate scaling parameters to link the remaining scores to the YSR metric at age 14. To do so, we repeated Steps 1–3 by (a) linking the scores at age 16 to the newly scaled scores at age 15; (b) linking scores at age 17 to the newly scaled scores at age 16; and (c) and so forth, until scores at all ages, including the YASR scores, had been linked to the target YSR

scale at age 14. The scaling parameters include an intercept parameter,  $B$ , and a slope parameter,  $A$ , that link the trait level scores at one age to the trait level scores at the prior age, by linking the discrimination and severity parameters at the two ages using the following formulas (Kolen & Brennan, 2014):

$$a(\text{age}_i) = \frac{a(\text{age}_j)}{A} \quad (1)$$

$$b(\text{age}_i) = A \times b(\text{age}_j) + B \quad (2)$$

where  $a(\text{age}_i)$  and  $a(\text{age}_j)$  represent the discrimination parameter for the common items at age  $i$  and age  $j$ , respectively;  $b(\text{age}_i)$  and  $b(\text{age}_j)$  represent the severity parameter for the common items at age  $i$  and age  $j$ , respectively;  $A$  represents the slope scaling parameter, and  $B$  represents the intercept scaling parameter.

4. We then used the scaling parameters to calculate the trait level scores at each age on the same scale. We used expected a posteriori (EAP) factor scores as individuals' trait level scores of internalizing problems. The  $A$  and  $B$  transformation constants rescale the standard deviation and mean, respectively, of the trait level scores to put the measures on a comparable scale, while still retaining changes in means and variances over time (based on the changes in means and variances of the common items). The vertically scaled scores were calculated by the following formula (Kolen & Brennan, 2014):

$$\theta(\text{age}_{14}) = A \times \theta(\text{age}_j) + B \quad (3)$$

where  $\theta(\text{age}_{14})$  represents the vector of trait level scores (i.e., factor scores) on the metric of the YSR at age 14, and  $\theta(\text{age}_j)$  represents the vector of trait level scores on the YSR or YASR at the remaining ages. When linking and chaining were completed, all scores were placed on the YSR age 14 metric.

**Growth curve model.** After vertically scaling the scores of internalizing problems to be on the same scale, we then examined individuals' trajectories of internalizing problems. To examine individuals' growth curves of vertically scaled internalizing problems, we used the lme function of the nlme package (Pinheiro, Bates, Debra, & Sarkar, 2009) in R for hierarchical linear modeling (HLM).<sup>2</sup> HLM can handle missingness and unbalanced data (Singer & Willett, 2003). We compared linear and curvilinear (polynomial) forms of growth using nested model comparisons with likelihood ratio tests. After settling on a form of growth, we examined sex and ethnicity as predictors of individuals' trajectories of internalizing problems.

## Results

### IRT Models

After determining that we approximately met IRT assumptions (Appendix S1) and observed only modest differential item functioning (DIF; Appendix S2), we fit a separate IRT graded response model at each age using the self-reported questionnaire items of internalizing problems. RMSEA estimates ranged from .058 to .072, depending on the year. CFI estimates ranged from .92 to .97, depending on the year. Thus, our model fit was adequate to good.

### Linking the YASR (and YSR) Scores to the Scale of the YSR at Age 14

Next, we linked the YASR and YSR scores so that scores on the two measures were on the same scale and could be compared. To link the two measures, we rescaled scores at all ages to the scale of the YSR at age 14 (see Steps 1–4 from the Vertical Scaling section of the Statistical Analysis section of the Method section). First, we examined scores on the 17 common items (i.e., the items that were common to both the YSR and YASR; see Table 2 and Panel A of Figure 2). Second, we fit separate IRT models at each age (see previous section).

Third, we used vertical scaling to put the IRT scores on the same scale. We calculated linear scaling parameters (slope:  $A$ ; intercept:  $B$ ) that linked the IRT scores at each age to the scores at the preceding age. The linear scaling parameters are in Table 2.<sup>3</sup>

Fourth, we used the scaling parameters to calculate individuals' internalizing problem scores on the same scale as the YSR at age 14. Age 15 scores were rescaled to the target scale of the age 14 scores by multiplying the age 15 scores by 1.091 and adding 0.017. We then applied linking and chaining to link the remaining scores to the YSR metric at age 14. To do so, we repeated steps 1–4 by linking the scores at age 16 to the newly scaled scores at age 15, linking scores at age 17 to the newly scaled scores at age 16, and so forth. For instance, age 16 scores were rescaled to the target scale at age 14 by the scaling parameters at age 16 (to transfer the scores to the age 15 metric) and then by the scaling parameters at age 15 (to transfer the scores to the age 14 metric). We applied this process of linking and chaining until all scores, including YASR scores, had been rescaled to the metric of the YSR at age 14.<sup>4</sup>

The mean and standard deviation of the vertically scaled scores are in Table 2. Participants' mean internalizing problem scores, after rescaling the YASR scores to be on the same metric as the YSR, are depicted in Panel B of Figure 2. Notably, the scores retained a highly similar pattern of mean scores by age when examining the rescaled total scores compared with when examining just the common items (see Panel A of Figure 2). Thus, the IRT approach to vertical scaling successfully retained mean-level change when rescaling the YASR scores to be on the same metric as the YSR while still using a more comparable scale. Moreover, vertically scaled scores from IRT were highly correlated with vertically scaled scores from Thurstone scaling ( $r = .95-.97$ , depending on the year).

<sup>2</sup> Although we considered multiple imputation approaches to handle missingness, to fairly compare the approach of using the common items with using the rescaled scores, we used only the observed data.

<sup>3</sup> Scaling parameters where  $A$  equals 1 and  $B$  equals 0 would represent no adjustment, so greater deviations from those values reflect greater adjustment to put the scores on the same scale. Notably, all of the scaling parameters for scores at adjacent ages were relatively close to these values ( $A \approx 1$  and  $B \approx 0$ ), indicating that only small adjustments were necessary to link the scores at adjacent ages.

<sup>4</sup> IRT factor scores have a mean of 0 and a standard deviation of 1, so the IRT-based internalizing problem scores at age 14 have a relatively normal distribution with a mean of 0 and standard deviation of 1. Scores at subsequent ages were linked to the target scale at age 14, so deviations from a mean of 0 and a standard deviation of 1 reflect changes in means and variances over time. For an example calculation of linking and chaining, see the note of Table 2.



### Growth Curve Model

To examine growth curves, we first compared a linear growth curve model with polynomial forms of change in HLM to identify the best-fitting form of change for the rescaled internalizing problem scores. A model with random linear slopes fit better than a model with a linear slope component that was fixed across individuals (i.e., fixed linear slopes;  $\chi^2(2) = 486.49, p < .001$ ). A model with a random linear slope component and a fixed quadratic component fit better than a model with only random linear slopes,  $\chi^2(1) = 24.25, p < .001$ . A model with a random linear slope component and a random quadratic component fit better than a model with a random linear slope component and a fixed quadratic component,  $\chi^2(3) = 64.37, p < .001$ , and was the best fitting model (model fit did not significantly improve when adding a fixed cubic component:  $\chi^2(1) = 2.25, p = .133$ ). Thus, a quadratic form of change was the best-fitting form of change for the rescaled internalizing problem scores. Individuals' quadratic trajectories, and the average quadratic trajectory for males and females are

depicted in Figure 3. The average quadratic trajectory showed slight decreases over time, primarily for females.

Overall, the growth curves showed little curvature, which would be consistent with evidence that likelihood ratio tests may be sensitive to small fit differences with larger sample sizes (Tomarken & Waller, 2003). Thus, the polynomial growth terms may have overfit the data, especially given the lengthy developmental span. Moreover, there are difficulties in interpreting and replicating findings from polynomial growth models, and mapping polynomial growth terms onto developmental theory (Grimm, Ram, & Hamagami, 2011). For these reasons, for comparing the common items to the rescaled scores and for examining the predictors of change in internalizing problems, we examined the general form of change by examining the linear model for ease of interpretation.

In the linear growth curve model with no predictors of the intercepts or slopes, intercepts reflected an individual's estimated initial level of internalizing problems at age 14. Slopes reflected participants' linear change in internalizing problems over time.

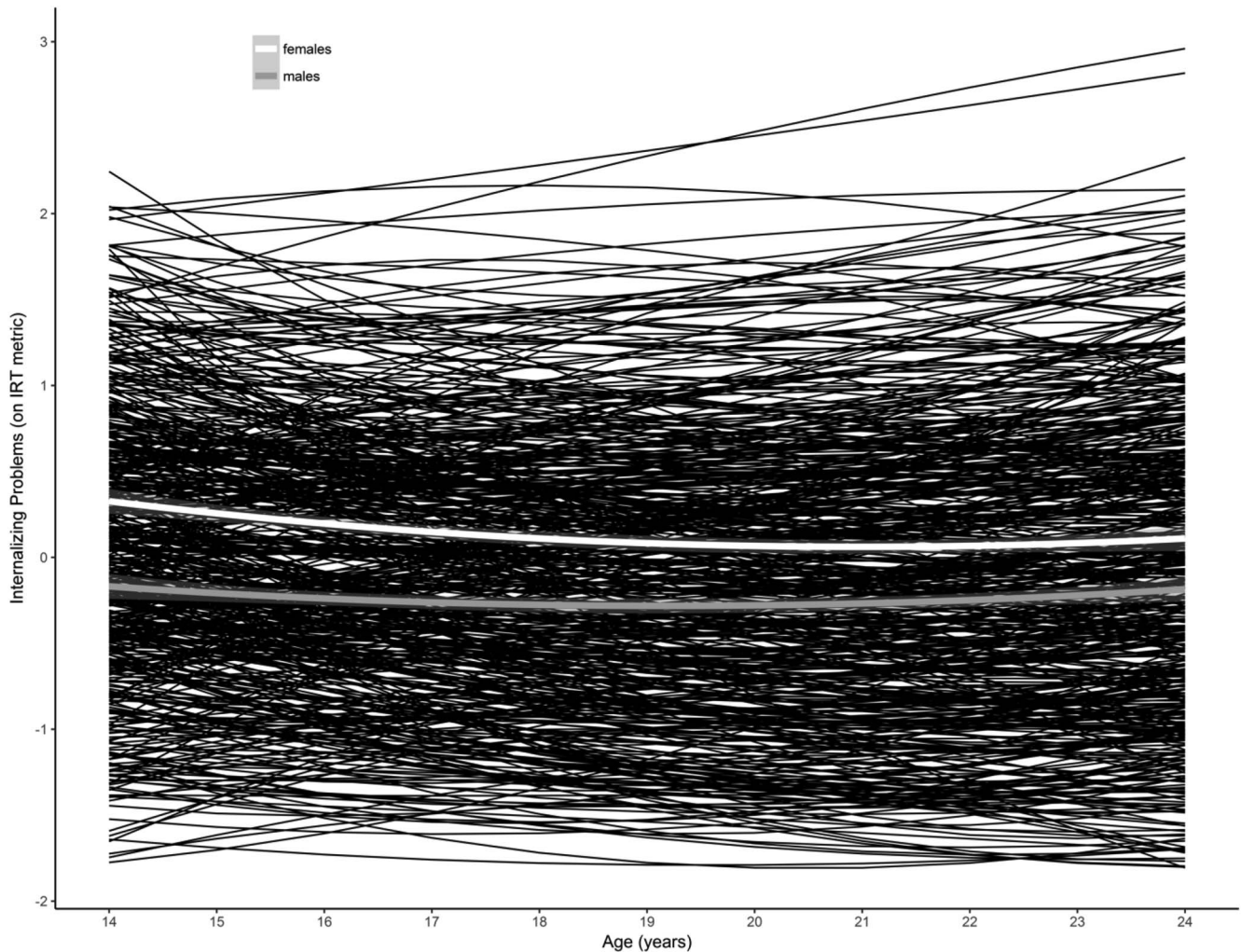


Figure 3. Individuals' fitted quadratic trajectories of internalizing problems in black (on IRT metric of YSR at age 14). Average quadratic trajectory for females in white. Average quadratic trajectory for males in gray.

There was a significant negative mean of the slopes ( $B = -0.01$ ,  $\beta = -0.04$ ,  $t(3980) = -2.21$ ,  $p = .027$ ). In a similar growth curve model examining the trajectories of scores on the common items, however, the mean of the slopes was not significant ( $B = -0.03$ ,  $\beta = -0.02$ ,  $t(3980) = -1.08$ ,  $p = .282$ ).

We conducted sensitivity analyses to determine the sensitivity of our findings to other vertical scaling approaches. The factor scores from a partially constrained multiple group (concurrent calibration) IRT model (within-item parameter constraints across time for non-DIF parameters) showed similar evidence of a negative mean of the slopes ( $B = -0.009$ ,  $\beta = -0.03$ ,  $t(3980) = -2.00$ ,  $p = .045$ ). There was also a negative mean of the slopes using factor scores from a comparable IRT model that excluded the items showing DIF ( $B = -0.01$ ,  $\beta = -0.04$ ,  $t(3980) = -2.16$ ,  $p = .031$ ). The Thurstone approach to vertical scaling also showed evidence of a negative mean of the slopes ( $B = -0.08$ ,  $\beta = -0.03$ ,  $t(3980) = -1.94$ ,  $p = .053$ ), at a trend-level.

In addition to differences in the form of change for the age-relevant items versus common items at the *group-level*, there were also differences at the *individual-level*. Some participants showed *decreases* in internalizing problems over time when using the age-relevant items while they showed *increases* in internalizing problems when using the common items (or vice versa). The participants who showed decreases using the age-relevant items and increases using the common items presumably had higher levels of internalizing problems on the *noncommon* items of the YSR (i.e., items that were on the Internalizing scale of the YSR but not the YASR) or lower levels on the *noncommon* items of the YASR (compared with the other participants). Because the somatic complaints subscale was included in the internalizing scale of the YSR but not YASR, the majority (nine items, 60%) of the non-common internalizing items of the YSR were items assessing somatic complaints. Therefore, we examined participants' levels of somatic complaints on the YSR. Consistent with expectations, participants who showed decreases in internalizing problems using the age-relevant items but increases using the common items

showed higher mean levels of somatic complaints from ages 14–19 ( $M = 3.36$ ) than participants who did not ( $M = 1.81$ ;  $t(29.07) = -3.69$ ,  $p < .001$ ). The reverse was also true; participants who showed increases in internalizing problems using the age-relevant items but decreases using the common items, showed lower mean levels of somatic complaints from ages 14–19 ( $M = 0.66$ ) than participants who did not ( $M = 1.95$ ;  $t(31.93) = 6.11$ ,  $p < .001$ ).

We then examined sex and ethnicity as predictors of the intercepts and linear slopes of the rescaled internalizing problem scores (see Table 3). The mean of the linear slopes was not significant when controlling for the other model predictors. Females showed higher intercepts than males, and showed a trend toward greater decreases over time compared to males. Although African Americans showed a trend toward lower intercepts, African Americans and those of “other” ethnicity did not significantly differ from European Americans in their intercepts or linear slopes.

## Discussion

Heterotypic continuity, the change in the manifestation of a construct or process over time, presents challenges to studying individuals over lengthy spans of development, and may necessitate using different measures over time. We examined self-reports of internalizing problems on the YSR from ages 14 to 19 and on the YASR from ages 20 to 24. The YSR internalizing scale includes items reflecting anxiety, depression, and somatic complaints, whereas its YASR counterpart includes items reflecting anxiety and depression but not somatic complaints. The challenge is measuring actual change rather than change in the meaning of the measures. We applied a vertical scaling technique to account for the heterotypic continuity of internalizing problems and the change in measurement.

Applying vertical scaling with age-relevant items to account for the heterotypic continuity of internalizing problems, we observed a pattern of means by age at the group-level that was similar to what we would have observed had we used the common items (see

Table 3  
*Linear Growth Curve Model of Internalizing Problems*

Variable	<i>B</i>	$\beta$	<i>SE</i>	<i>df</i>	<i>p</i>
Intercept	-.197	.013	.060	3977	.001
Time	.000	-.036	.008	3977	.963
Predictors of the intercepts					
Female	.477	.183	.081	539	<.001
African American	-.189	-.089	.112	539	.091
Other ethnicity	-.329	-.024	.343	539	.337
Predictors of the slopes					
Female	-.018	-.029	.010	3977	.078
African American	-.014	-.016	.015	3977	.329
Other ethnicity	.022	.009	.043	3977	.601
Variance components					
	<i>SD</i>				
Intercept	.83				
Time	.10				
Residual	.58				
Correlation between intercept and slope	$r = -.41$				
Model Pseudo- $R^2$	.753				

*Note.* The model's pseudo  $R^2$  was calculated as the squared correlation between the model's fitted and observed values (Singer & Willett, 2003).

Figure 2), but with important differences. The age-relevant items showed a *group-level* pattern of means by age similar to the results with the common items, but the age-relevant items resulted in more construct-valid scores of internalizing problems at the *individual-level*. The age trends of the mean values of the *observed* scores differed from the means of individuals' slopes based on *model-fitted* values in growth curve models, which fit lines through all available time points, and essentially interpolate missing values based on the individual's other time points and the other individuals' trajectories (i.e., shrinkage). We found that vertical scaling made small adjustments to the scores (see Table 2), but these subtle adjustments resulted in potentially meaningful differences in the individuals' and group trajectories. Because the common items ignored the age-specific manifestations of internalizing problems, for example, somatic complaints, some participants showed *decreases* in internalizing problems when we used their ratings on the age-relevant items but *increases* when we used only the common items, and other participants showed the opposite pattern. The differences in individuals' trajectories using the age-relevant items versus the common items could explain differences we observed in the group-level trajectories using the age-relevant items versus the common items. Although prior research is mixed, some studies have shown decreases in the prevalence of internalizing disorders from adolescence to adulthood (Costello, Cope-land, & Angold, 2011). We observed group-level decreases in internalizing problems from adolescence to early adulthood using the construct-valid items. Using only the common items, however, we observed no significant change in internalizing problems over time. Discarding items (e.g., somatic complaints) that were relevant to internalizing problems during some developmental periods but not other developmental periods (i.e., using only the common items) resulted in a loss of information that may have made the measure less sensitive to developmental change. Thus, accounting for heterotypic continuity could have theoretical and practical advantages over ignoring heterotypic continuity by using only the common items across time. Future research should further examine the potential reasons why the approaches may differ in their developmental inferences.

Accounting for heterotypic continuity allowed us to examine predictors of individuals' trajectories over a lengthy developmental span. We observed that females showed higher levels of internalizing problems than males at age 14, and there was a trend toward females showing greater decreases over time compared with males. As shown in Figure 2, we found the greatest difference between females' and males' levels of internalizing problems around ages 15–18, which is consistent with Hankin et al. (1998), and the greatest level of internalizing problems around age 15, consistent with Adkins et al. (2009). We also found a trend toward lower levels of internalizing problems among African Americans compared with European Americans.

Despite evidence of several items showing modest DIF over time, the overall theoretical and empirical evidence suggests we measured the same construct in an equivalent way across time. Although longitudinal measurement invariance should be tested, establishing strict longitudinal measurement invariance is unnecessary in the case of heterotypic continuity because the meaning of the measures is expected to change with changes in the manifestation of the construct (Petersen et al., 2016). Research has demonstrated that models with failed longitudinal measurement invari-

ance can yield valid inferences in the context of heterotypic continuity (Edwards & Wirth, 2012). Removing items/measures that show DIF or failed measurement invariance over time is not necessarily recommended in the case of heterotypic continuity (Knight & Zerr, 2010). Removing items or measures can result in a less representative sample of the content of the construct (i.e., lower content validity), and some items might be expected to change in their discrimination or severity over time given heterotypic continuity, and yet remain construct-valid. Discarding them would be removing important and meaningful developmental information about the construct. Discarding construct-valid items showing DIF or failed measurement invariance would be akin to using only the common items, which we argue is highly problematic (and violates content validity). Nevertheless, we observed similar results when we excluded DIF items, suggesting that DIF did not compromise the findings.

In addition to empirical considerations, there are important theoretical considerations regarding whether one is measuring the same construct across time in an equivalent way (construct validity invariance). First, the Achenbach scales are widely used measures of internalizing problems; they were derived empirically, and have strong validity, including content validity, construct validity, and criterion-related validity (Sattler & Hoge, 2006). Second, the items were selected based on theory and on the known heterotypic continuity of internalizing problems—we used the age-relevant items instead of discarding construct-relevant items that were not present in both forms of the measure. Third, we observed strong internal consistency of the items at each age, and the items showed strong cross-time continuity (see Table 1). Fourth, the items showed convergent and discriminant validity with respect to externalizing problems. Fifth, the trajectories showed construct validity: their pattern was similar with prior findings. Finally, we observed similar trajectories with multiple approaches to vertical scaling, including separate IRT estimation with linking, concurrent calibration in IRT, and Thurstone scaling. Thus, we feel there is strong theoretical and empirical evidence for using the internalizing scales of the YSR/YASR as they are constructed for measuring the same construct of internalizing problems in an equivalent way over time in the present study.

### Alternative Approaches to Vertical Scaling

We applied the widely used IRT approach to vertical scaling, which uses a latent variable approach. There are alternative approaches to vertical scaling. Thurstone scaling, an observed score approach, may be more practical than IRT in some situations for vertical scaling. First, IRT requires large sample sizes for accurate estimation. Second, IRT generally requires dichotomous, polytomous, or categorical items instead of continuous measures (unless moving to a SEM framework). Third, except for advanced and cutting-edge multidimensional IRT techniques, most IRT applications require items that are unidimensional. These requirements pose challenges for psychological constructs, which are often multifaceted and measured using various metrics. Nevertheless, (unidimensional) IRT is often employed for vertical scaling, and the findings are often consistent with Thurstone scaling (Becker & Forsyth, 1992), as they were in the present study.



## Implications for Developmental Psychology

We are unaware of other studies that have examined individuals' trajectories of broadband internalizing problems from adolescence to adulthood. The present results show how internalizing problems developed across an important developmental transition. Broadband internalizing problems peaked in mid-to-late adolescence and decreased into adulthood, similar to patterns shown for depression (Adkins et al., 2009). Further, the decreases in internalizing problems were detected after we accounted for their heterotypic continuity using vertical scaling. The findings of the present study are novel, but the statistical approach is not. Previous research has (a) used vertical scaling to link different measures on the same comparable scale (Kolen & Brennan, 2014); (b) measured change with different measures (McArdle, Grimm, Hamagami, Bowles, & Meredith, 2009); and (c) used changing items to account for the heterotypic continuity of psychopathology based on theory (Petersen, Bates, Dodge, Lansford, & Pettit, 2015). What is especially novel in the present study is the assembling of these techniques to demonstrate how to use vertical scaling and changing measures to account for heterotypic continuity and measure individuals' change in constructs showing heterotypic continuity. We feel this is a crucial theoretical and empirical advance, especially for developmental theory. The vast majority of research in developmental psychology has examined trajectories using the same measures over time, which is a common practice with some advantages for model building, but which, as we argue next is often highly problematic for developmental theory.

When developmental psychologists have examined individuals' change in a construct using the same measures over time, in the traditional way, using either all available items or only age-common items, this creates a theoretical and empirical problem when the construct shows heterotypic continuity, that is, change in its manifestation over time. Using all available items over time violates construct validity because, to one degree or another, the same items do not consistently reflect the same construct over time. Using only age-common items violates content validity because the measure is not assessing the construct as a whole, including its age-specific manifestations. Moreover, not only are there theoretical reasons to use different measures over time in the context of heterotypic continuity, there are likely empirical advantages of using different measures over time, as well. We showed that using different measures (i.e., all construct-valid items) over time may be more sensitive to developmental change than using only age-common items.

## Strengths and Limitations

The present study had key strengths. First, it examined the development of broadband internalizing problems over a lengthy span of development across the important developmental transition from adolescence to adulthood. Second, it accounted for the heterotypic continuity of internalizing problems when examining individuals' trajectories, and compared the approach with traditional approaches that ignore heterotypic continuity. Third, it examined the form of change of internalizing problems and sex and ethnicity as predictors of the trajectories. Fourth, it considered multiple approaches to vertical scaling, each with different as-

sumptions, and found substantially similar results with each method, providing greater confidence in the findings.

The study also had limitations. We did not examine trajectories of individual items or subdimensions of internalizing problems (e.g., anxiety or depression). One can always reduce to a lower level subunit, however. Internalizing problems have an empirically derived factor structure, so we believe there is theoretical reason for this level of analysis. In addition, the subdimensions of anxiety and depression, themselves, like most behavior trait measures, are heterogeneous and involve behaviors whose meaning would depend on age, so they would likely demonstrate heterotypic continuity, as well.

## Conclusion

The present study applied vertical scaling to account for the heterotypic continuity of internalizing problems from adolescence to adulthood. Vertical scaling allowed us to place scores from two measures on the same scale. Accounting for heterotypic continuity by using all developmentally relevant items may have been more sensitive to developmental change in internalizing problems than was ignoring heterotypic continuity by using the same items across major stages of development. Using vertical scaling, internalizing problems peaked in mid-to-late adolescence and decreased into adulthood. Vertical scaling may be a useful approach to measuring individuals' developmental trajectories in constructs that change in their manifestation over time.

## References

- Achenbach, T. M. (1991). *Manual for the youth self-report and 1991 profile*. Burlington, VT: University of Vermont, Department of Psychiatry.
- Achenbach, T. M. (1997). *Manual for the young adult self-report and young adult behavior checklist*. Burlington, VT: University of Vermont, Department of Psychiatry.
- Achenbach, T. M., & Edelbrock, C. S. (1978). The classification of child psychopathology: A review and analysis of empirical efforts. *Psychological Bulletin*, *85*, 1275–1301. <http://dx.doi.org/10.1037/0033-2909.85.6.1275>
- Adkins, D. E., Wang, V., & Elder, G. H. (2009). Structure and stress: Trajectories of depressive symptoms across adolescence and young adulthood. *Social Forces*, *88*, 31–60. <http://dx.doi.org/10.1353/sof.0.0238>
- Avenevoli, S., & Steinberg, L. (2002). The continuity of depression across the adolescent transition. In W. R. Hayne & K. Robert (Eds.), *Advances in child development and behavior* (Vol. 28, pp. 139–173). San Diego, CA: Academic Press. [http://dx.doi.org/10.1016/S0065-2407\(02\)80064-7](http://dx.doi.org/10.1016/S0065-2407(02)80064-7)
- Becker, D. F., & Forsyth, R. A. (1992). An empirical investigation of Thurstone and IRT methods of scaling achievement tests. *Journal of Educational Measurement*, *29*, 341–354. <http://dx.doi.org/10.1111/j.1745-3984.1992.tb00382.x>
- Betts, K. S., Baker, P., Alati, R., McIntosh, J. E., Macdonald, J. A., Letcher, P., & Olsson, C. A. (2016). The natural history of internalizing behaviours from adolescence to emerging adulthood: Findings from the Australian Temperament Project. *Psychological Medicine*, *46*, 2815–2827. <http://dx.doi.org/10.1017/S0033291716001495>
- Broeren, S., Muris, P., Diamantopoulou, S., & Baker, J. R. (2013). The course of childhood anxiety symptoms: Developmental trajectories and child-related factors in normal children. *Journal of Abnormal Child Psychology*, *41*, 81–95. <http://dx.doi.org/10.1007/s10802-012-9669-9>



- Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, *48*, 1–29. <http://dx.doi.org/10.18637/jss.v048.i06>
- Chen, F. R., & Jaffee, S. R. (2015). The heterogeneity in the development of homotypic and heterotypic antisocial behavior. *Journal of Developmental and Life-Course Criminology*, *1*, 269–288. <http://dx.doi.org/10.1007/s40865-015-0012-3>
- Chen, W.-H., & Thissen, D. (1997). Local dependence indexes for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, *22*, 265–289. <http://dx.doi.org/10.2307/1165285>
- Costello, E. J., Copeland, W., & Angold, A. (2011). Trends in psychopathology across the adolescent years: What changes when children become adolescents, and when adolescents become adults? *Journal of Child Psychology and Psychiatry*, *52*, 1015–1025. <http://dx.doi.org/10.1111/j.1469-7610.2011.02446.x>
- Côté, S. M., Boivin, M., Liu, X., Nagin, D. S., Zoccolillo, M., & Tremblay, R. E. (2009). Depression and anxiety symptoms: Onset, developmental course and risk factors during early childhood. *Journal of Child Psychology and Psychiatry*, *50*, 1201–1208. <http://dx.doi.org/10.1111/j.1469-7610.2009.02099.x>
- Crocetti, E., Klimstra, T., Keijsers, L., Hale, W. W., III, & Meeus, W. (2009). Anxiety trajectories and identity development in adolescence: A five-wave longitudinal study. *Journal of Youth and Adolescence*, *38*, 839–849. <http://dx.doi.org/10.1007/s10964-008-9302-y>
- De Los Reyes, A., & Kazdin, A. E. (2005). Informant discrepancies in the assessment of childhood psychopathology: A critical review, theoretical framework, and recommendations for further study. *Psychological Bulletin*, *131*, 483–509. <http://dx.doi.org/10.1037/0033-2909.131.4.483>
- Dodge, K. A., Bates, J. E., & Pettit, G. S. (1990). Mechanisms in the cycle of violence. *Science*, *250*, 1678–1683. <http://dx.doi.org/10.1126/science.2270481>
- Eaton, N. R., Krueger, R. F., Markon, K. E., Keyes, K. M., Skodol, A. E., Wall, M., . . . Grant, B. F. (2013). The structure and predictive validity of the internalizing disorders. *Journal of Abnormal Psychology*, *122*, 86–92. <http://dx.doi.org/10.1037/a0029598>
- Edwards, M. C., & Wirth, R. J. (2012). Valid measurement without factorial invariance: A longitudinal example. In G. R. Hancock & J. R. Harring (Eds.), *Advances in longitudinal methods in the social and behavioral sciences* (pp. 289–311). Charlotte, NC: Information Age Publishing.
- Fanti, K. A., & Henrich, C. C. (2010). Trajectories of pure and co-occurring internalizing and externalizing problems from age 2 to age 12: Findings from the National Institute of Child Health and Human Development Study of Early Child Care. *Developmental Psychology*, *46*, 1159–1175. <http://dx.doi.org/10.1037/a0020659>
- Fennessy, L. M. (1995). *The impact of local dependencies on various IRT outcomes* (Doctoral dissertation). Available from ProQuest Dissertations & Theses database. (UMI No. 9524701)
- Gilliom, M., & Shaw, D. S. (2004). Codevelopment of externalizing and internalizing problems in early childhood. *Development and Psychopathology*, *16*, 313–333. <http://dx.doi.org/10.1017/S0954579404044530>
- Grimm, K. J., Ram, N., & Hamagami, F. (2011). Nonlinear growth curves in developmental research. *Child Development*, *82*, 1357–1371. <http://dx.doi.org/10.1111/j.1467-8624.2011.01630.x>
- Hale, W. W., III, Raaijmakers, Q., Muris, P., van Hoof, A., & Meeus, W. (2008). Developmental trajectories of adolescent anxiety disorder symptoms: A 5-year prospective community study. *Journal of the American Academy of Child & Adolescent Psychiatry*, *47*, 556–564. <http://dx.doi.org/10.1097/CHI.0b013e3181676583>
- Hankin, B. L., Abramson, L. Y., Moffitt, T. E., Silva, P. A., McGee, R., & Angell, K. E. (1998). Development of depression from preadolescence to young adulthood: Emerging gender differences in a 10-year longitudinal study. *Journal of Abnormal Psychology*, *107*, 128–140. <http://dx.doi.org/10.1037/0021-843X.107.1.128>
- Knight, G. P., & Zerr, A. A. (2010). Informed theory and measurement equivalence in child development research. *Child Development Perspectives*, *4*, 25–30. <http://dx.doi.org/10.1111/j.1750-8606.2009.00112.x>
- Kolen, M. J., & Brennan, R. L. (2014). *Test equating, scaling, and linking: Methods and practices* (3rd ed.). New York, NY: Springer. <http://dx.doi.org/10.1007/978-1-4939-0317-7>
- Leadbeater, B., Thompson, K., & Gruppiso, V. (2012). Co-occurring trajectories of symptoms of anxiety, depression, and oppositional defiance from adolescence to young adulthood. *Journal of Clinical Child and Adolescent Psychology*, *41*, 719–730. <http://dx.doi.org/10.1080/15374416.2012.694608>
- Letcher, P., Sanson, A., Smart, D., & Toumbourou, J. W. (2012). Precursors and correlates of anxiety trajectories from late childhood to late adolescence. *Journal of Clinical Child and Adolescent Psychology*, *41*, 417–432. <http://dx.doi.org/10.1080/15374416.2012.680189>
- Letcher, P., Smart, D., Sanson, A., & Toumbourou, J. W. (2009). Psychosocial precursors and correlates of differing internalizing trajectories from 3 to 15 years. *Social Development*, *18*, 618–646. <http://dx.doi.org/10.1111/j.1467-9507.2008.00500.x>
- Markon, K. E., Chmielewski, M., & Miller, C. J. (2011). The reliability and validity of discrete and continuous measures of psychopathology: A quantitative review. *Psychological Bulletin*, *137*, 856–879. <http://dx.doi.org/10.1037/a0023678>
- Mathiesen, K. S., Sanson, A., Stoolmiller, M., & Karevold, E. (2009). The nature and predictors of undercontrolled and internalizing problem trajectories across early childhood. *Journal of Abnormal Child Psychology*, *37*, 209–222. <http://dx.doi.org/10.1007/s10802-008-9268-y>
- McArdle, J. J., Grimm, K. J., Hamagami, F., Bowles, R. P., & Meredith, W. (2009). Modeling life-span growth curves of cognition using longitudinal data with multiple samples and changing scales of measurement. *Psychological Methods*, *14*, 126–149. <http://dx.doi.org/10.1037/a0015857>
- Meade, A. W. (2010). A taxonomy of effect size measures for the differential functioning of items and scales. *Journal of Applied Psychology*, *95*, 728–743. <http://dx.doi.org/10.1037/a0018966>
- Meadows, S. O., Brown, J. S., & Elder, G. H., Jr. (2006). Depressive symptoms, stress, and support: Gendered trajectories from adolescence to young adulthood. *Journal of Youth and Adolescence*, *35*, 89–99. <http://dx.doi.org/10.1007/s10964-005-9021-6>
- Miers, A. C., Blöte, A. W., de Rooij, M., Bokhorst, C. L., & Westenberg, P. M. (2013). Trajectories of social anxiety during adolescence and relations with cognition, social competence, and temperament. *Journal of Abnormal Child Psychology*, *41*, 97–110. <http://dx.doi.org/10.1007/s10802-012-9651-6>
- Moeller, J. (2015). A word on standardization in longitudinal studies: Don't. *Frontiers in Psychology*, *6*, 1389. <http://dx.doi.org/10.3389/fpsyg.2015.01389>
- Morin, A. J. S., Maïano, C., Nagengast, B., Marsh, H. W., Morizot, J., & Janosz, M. (2011). General growth mixture analysis of adolescents' developmental trajectories of anxiety: The impact of untested invariance assumptions on substantive interpretations. *Structural Equation Modeling*, *18*, 613–648. <http://dx.doi.org/10.1080/10705511.2011.607714>
- Morizot, J., Ainsworth, A. T., & Reise, S. P. (2007). Toward modern psychometrics: Application of item response theory models in personality research. In R. W. Robins, R. C. Fraley, & R. F. Krueger (Eds.), *Handbook of research methods in personality psychology* (pp. 407–421). New York, NY: Guilford Press.
- Petersen, I. T., Bates, J. E., Dodge, K. A., Lansford, J. E., & Pettit, G. S. (2015). Describing and predicting developmental profiles of externalizing problems from childhood to adulthood. *Development and Psychopathology*, *27*, 791–818. <http://dx.doi.org/10.1017/S0954579414000789>
- Petersen, I. T., Hoyniak, C. P., McQuillan, M. E., Bates, J. E., & Staples, A. D. (2016). Measuring the development of inhibitory control: The

- challenge of heterotypic continuity. *Developmental Review*, 40, 25–71. <http://dx.doi.org/10.1016/j.dr.2016.02.001>
- Pinheiro, J., & Bates, D., DebRoy, S., & Sarkar, D., & the R Core team. (2009). *nlme: Linear and nonlinear mixed effects models*. R package version 3.1–93. Retrieved from <http://cran.r-project.org/web/packages/nlme/index.html>
- Raju, N. S. (1990). Determining the significance of estimated signed and unsigned areas between two item response functions. *Applied Psychological Measurement*, 14, 197–207. <http://dx.doi.org/10.1177/014662169001400208>
- Ryan, N. D., Puig-Antich, J., Ambrosini, P., Rabinovich, H., Robinson, D., Nelson, B., . . . Twomey, J. (1987). The clinical picture of major depression in children and adolescents. *Archives of General Psychiatry*, 44, 854–861. <http://dx.doi.org/10.1001/archpsyc.1987.01800220016003>
- Sattler, J. M., & Hoge, R. D. (2006). *Assessment of children: Behavioral, social, and clinical foundations* (5th ed.). San Diego, CA: Jerome M. Sattler, Publisher, Inc.
- Singer, J. D., & Willett, J. B. (2003). *Applied longitudinal data analysis: Modeling change and event occurrence*. New York, NY: Oxford University Press, Inc. <http://dx.doi.org/10.1093/acprof:oso/9780195152968.001.0001>
- Snyder, H. R., Young, J. F., & Hankin, B. L. (2017). Strong homotypic continuity in common psychopathology-, internalizing-, and externalizing-specific factors over time in adolescents. *Clinical Psychological Science*, 5, 98–110. <http://dx.doi.org/10.1177/2167702616651076>
- Sterba, S. K., Prinstein, M. J., & Cox, M. J. (2007). Trajectories of internalizing problems across childhood: Heterogeneity, external validity, and gender differences. *Development and Psychopathology*, 19, 345–366. <http://dx.doi.org/10.1017/S0954579407070174>
- Tomarken, A. J., & Waller, N. G. (2003). Potential problems with “well fitting” models. *Journal of Abnormal Psychology*, 112, 578–598. <http://dx.doi.org/10.1037/0021-843X.112.4.578>
- Toumbourou, J. W., Williams, I., Letcher, P., Sanson, A., & Smart, D. (2011). Developmental trajectories of internalising behaviour in the prediction of adolescent depressive symptoms. *Australian Journal of Psychology*, 63, 214–223. <http://dx.doi.org/10.1111/j.1742-9536.2011.00023.x>
- van der Ark, L. A. (2007). Mokken scale analysis in R. *Journal of Statistical Software*, 20, 19.
- Weeks, J. P. (2010). plink: An R package for linking mixed-format tests using IRT-based methods. *Journal of Statistical Software*, 35, 33.

Received November 12, 2016

Revision received August 21, 2017

Accepted August 30, 2017 ■

### Correction to Bratt et al. (2018)

In the article “Perceived Age Discrimination Across Age in Europe: From an Ageing Society to a Society for All Ages,” by Christopher Bratt, Dominic Abrams, Hannah J. Swift, Christin-Melanie Vauclair, and Sibila Marques (*Developmental Psychology*, 2018, Vol. 54, No. 1, pp. 167–180. <http://dx.doi.org/10.1037/dev0000398>), the copyright license has been changed to the Creative Commons CC-BY Attribution License. The online version of this article has been corrected.

<http://dx.doi.org/10.1037/dev0000540>

## Supplementary Appendix

- Appendix S1. Item response theory (IRT) assumptions.
- Appendix S2. Differential item functioning (DIF).
- Appendix S3. Thurstone scaling.
- Table S1. Raw score frequency distributions on common items for females.
- Table S2. Raw score frequency distributions on common items for males.
- Table S3. Percentile ranks (divided by 100) on common items for females.
- Table S4. Percentile ranks (divided by 100) on common items for males.
- Table S5. Z-scores of common items for females.
- Table S6. Z-scores of common items for males.
- Table S7. Raw score frequency distributions on all items for females.
- Table S8. Raw score frequency distributions on all items for males.
- Table S9. Percentile ranks (divided by 100) on all items for females.
- Table S10. Percentile ranks (divided by 100) on all items for males.
- Table S11. Z-scores of all items for females.
- Table S12. Z-scores of all items for males.
- Table S13. Thurstone-scaled conversion table of YASR to YSR equivalents.
- Table S14. Linear growth curve model of Thurstone-scaled internalizing problems.
- Figure S1. Depiction of steps in vertical scaling using Thurstone scaling with a common-item design.
- Figure S2. Means of common items and all age-relevant Thurstone-scaled items over time.
- Figure S3. Individuals' fitted trajectories of Thurstone-scaled internalizing problems.

## Appendix S1

### Item Response Theory (IRT) assumptions

#### Method

We evaluated three IRT assumptions: (1) uni-dimensionality—the items have one predominant dimension reflecting the underlying (latent) trait (i.e., internalizing problems); (2) local independence—the items are uncorrelated when controlling for the latent dimension; and (3) monotonicity—the probability of endorsing a higher level on an item increases as the person’s trait level of internalizing problems increases. Our criterion for uni-dimensionality was a ratio of first to second eigenvalues of  $\geq 3.0$  for an unrotated factor solution (Morizot, Ainsworth, & Reise, 2007). We evaluated the local independence of items by examining the  $X^2$  local dependence (LD) statistic between each pair of items after controlling for the latent dimension (Chen & Thissen, 1997). We compared the LD statistics against a chi-square distribution using the mirt package (Chalmers, 2012) in R 3.3.2 to determine the proportion of pairwise items that showed greater dependency than would be expected by chance. We evaluated monotonicity by fitting non-parametric IRT models using Mokken scale analysis in the mokken package (van der Ark, 2007) in R. We examined the number of violations of monotonicity (i.e., decreases in the item step response function by rest score group) whose size was significantly greater than zero.

#### Results

In terms of dimensionality, the ratios of the first to second eigenvalues from unrotated factor solutions ranged from 3.22 to 3.90 across ages 14–24, suggesting that the data were “uni-dimensional enough” for IRT. In terms of local dependency, with an alpha level of .05 that was not corrected for multiple testing (4,650 pairwise associations), we observed that 4.7–8.8% of



pairwise associations had statistically significant dependency, depending on the year. Thus, after accounting for the expected Type I error rate of 5%, there were between 0–3.8% of associations showing greater linear dependency than would be expected by chance, depending on the year. Thus, there was modest evidence of some local dependency at some ages. Nevertheless, IRT is robust to low and moderate violations of the local independence assumption (Fennessy, 1995). In terms of monotonicity, no items showed statistically significant non-monotonicity at any ages.

### **Discussion**

Given evidence supporting that we approximately met the assumptions of IRT, we proceeded with the IRT approach to vertical scaling.

## Appendix S2

### Differential Item Functioning (DIF)

#### Method

After fitting IRT models, we examined whether there was differential item functioning (DIF) across age (comparable to tests of longitudinal measurement/factorial invariance). DIF examines whether the likelihood of endorsing a particular item differs between groups (in this case, between ages) for people with the same trait levels. DIF was explored using a multiple group framework in IRT in which item parameters were estimated simultaneously in the same model (i.e., concurrent calibration). In this framework, the baseline model was one in which item parameter estimates were allowed to vary *across* items, but item parameter estimates were constrained to be equal *within* item across time (allowing discrimination to differ from severity). To explore DIF, item parameters were iteratively allowed to vary across time. For example, the discrimination parameters were iteratively allowed to vary item by item to see if estimation of unique discrimination parameters at each age resulted in better model fit (based on nested model comparisons using chi-square difference tests). When exploring DIF for the discrimination parameters, a chi-square statistic with nine degrees of freedom was used to identify items with DIF. A similar procedure was used for the two severity terms, which resulted in a chi-square statistic with 18 degrees of freedom. To limit the impact of multiple testing ( $37 \text{ items} \times 2 \text{ parameters} = 74 \text{ tests}$ ), we set an alpha level of .01 for identifying DIF, resulting in cutoffs of chi-square statistics greater than 21.66 and 34.80 for discrimination and severity, respectively. We also examined DIF by sex.

We then examined the effect size of DIF. Using a framework first defined by Raju (1990) and discussed by Meade (2010), the signed and unsigned differences between expected

scores across ages were used to quantify magnitude of DIF. Both signed and unsigned differences were used to help explore whether DIF was non-uniform. With non-uniform DIF, the expected score curves across ages may cross and may cause the signed differences between expected scores to be zero, whereas the unsigned differences sum the absolute value of differences to approximate absolute expected score differences. Expected scores were estimated using a range of internalizing problem scores from eight standard deviations below the mean to eight standard deviations above the mean, and the average signed and unsigned difference were calculated for items showing DIF. Age 14 was used as the reference age for all calculations. The metric of these measures is in the raw score metric; for example, an effect size of 0.1 would indicate that scores at the focal age are 0.1 points larger compared to scores at age 14 (Meade, 2010).

Upon identifying an item as having parameters that differed across time, two additional IRT models were explored to assess the impact of DIF on the resulting internalizing problem scores. First, a partially constrained model was used that (a) constrained parameter estimates to equality within item across time for items showing no evidence of DIF and (b) freely estimated DIF parameters across time for items with evidence of DIF by age. Second, a model was explored that freely estimated all item parameters across ages. Both models were estimated using multiple group (concurrent calibration) IRT models with the mirt package in R. Model fit indices were used to identify which model was best fitting. Factor scores were generated to compare to the factor scores from IRT models fit separately by age, and were transformed to be on a common age 14 metric using vertical scaling (using calculations described in the Method section of the manuscript). Comparisons of model fit and factor scores were conducted as a sensitivity analysis to determine the sensitivity of the model (separate linking versus multiple

group; impact of DIF) on the internalizing problems factor scores.

## Results

We examined DIF by sex and by age. We are hesitant to interpret findings of DIF by sex because we had relatively small subgroups for fitting multiple group models by sex, which may have resulted in unreliable parameter estimates. We observed some instances of DIF by sex at some ages. Consistent with the interpretation that multiple group models by sex may have yielded unreliable parameter estimates, however, the items showing DIF were not consistent across ages, suggesting that items mostly did not reliably differ between males and females.

Initial exploration of DIF by age revealed that, out of 37 items, four items showed evidence of DIF in terms of discrimination and nine items showed evidence of DIF in terms of severity (all of which were common items). Three of these items showed evidence of DIF with respect to both discrimination and severity. No consistent trend was found with respect to the directionality of DIF. Some items showed increases in severity or discrimination with age, whereas other items showed decreases.

We examined the effect size of DIF. On average, DIF across time tended to have a small effect size, ranging from 0–0.16 raw score points. The signed and unsigned metrics were similar within an item across age, suggesting uniform DIF (i.e. the expected score curves did not cross), which is unsurprising because there was less evidence of DIF with respect to discrimination compared to severity.

The two multiple group models, one that allowed item parameters with evidence of DIF to be estimated freely across time and another that freely estimated all parameters, showed similar model fit. When using the chi-square model fit statistic, the model with all parameters freely estimated fit the data the best, which is unsurprising given the sample size and lengthy



developmental span. However, when accounting for model complexity with a statistic such as the Akaike Information Criterion (AIC), the partially constrained model showed evidence of best model fit. Thus, model fit was not *considerably* worse fitting (accounting for model complexity) when constraining the non-DIF parameters across time to help anchor the latent variable to the same scale, suggesting that we were measuring the same construct in an equivalent way over time. At the same time, the model with all parameter estimates freely estimated was the best overall fitting model, which supports our decision to use separate IRT estimation and linking in the context of heterotypic continuity over a lengthy developmental span.

## **Discussion**

We observed several items showing potential changes in discrimination or severity over time. Changes in severity are expected across a lengthy developmental span and are less likely than changes in discrimination to be serious threats to measuring the same construct. Compared to changes in severity, changes in discrimination are potentially more serious because they may reflect that an item does not tap the same construct at some ages. However, changes in discrimination may instead reflect meaningful developmental shifts in the construct (heterotypic continuity) when the items tap the theoretical content of the construct, as was likely the case in the present study given the strong empirical basis and content validity of the measure we used. Nevertheless, most of the items showing evidence of DIF showed changes in severity rather than discrimination, and effect sizes of DIF were small. Moreover, even for those items that changed in discrimination, they were still highly discriminating across time, further supporting that we were measuring the same construct at all ages. Despite considerable research on DIF and measurement invariance, there is not clear guidance in the literature on how to proceed in the case of DIF (or failed measurement invariance) because there is no test to determine whether the

difference reflects a change in the manifestation of the construct (i.e., heterotypic continuity), changes in the functioning of the measures, or some combination of the two (Knight & Zerr, 2010). Nevertheless, we examined the effect size of DIF. All instances of DIF had small effect sizes. Our vertical scaling approach accounted for DIF by estimating a separate IRT model at each age, thus allowing items' parameters to change over time, and using scaling parameters to link the scores across age and "smooth out" the DIF at the construct-level. In sum, there are theoretical and empirical considerations when determining whether we measured the same construct in an equivalent way over time, and the totality of the evidence suggests that we did. Importantly, we found the same results with (a) a partially constrained model (within-item parameter constraints across time for non-DIF parameters), (b) a model excluding the items showing DIF, (c) separate IRT estimation and linking, and (d) Thurstone scaling, providing further confidence in our findings.

## Appendix S3

### Thurstone Scaling

When the different measures have common items over time, the Thurstone scaling approach to vertical scaling uses common items administered across ages to link the measures on the same scale by aligning their percentile scores based on a range of  $z$ -scores on the common items. This is based on the assumption that the two age groups to be linked have the same form of distribution (i.e., are normally distributed on the underlying trait within group), and that the groups' scores on the measures might differ in their mean and standard deviation.

#### Method

Vertical scaling involves placing two measures that assess a similar construct but differ in difficulty/severity on the same scale. Ideally, the two measures should have some items with the same contents to ensure scores on the measures can be linked (i.e., made comparable). In the present study, we used the Thurstone scaling approach to vertical scaling (as described in Kolen & Brennan, 2014) to transform scores on the YASR to the scale of the YSR (in addition to the IRT approach described in the manuscript). See Figure S1 for a depiction of the Thurstone scaling approach to vertical scaling.

The YSR and YASR have different but overlapping item content, so we needed to put them on the same scale. We applied Thurstone scaling that scales the scores across the different measures using the items that are in common across both measures (i.e., a common-item design, see Figure 1). Although the common items are used to determine the general form of change on the same scale, all developmentally relevant, construct-valid items are used to estimate each person's trait level on this scale. To link two measures using Thurstone scaling in a common-item design,  $z$ -scores are calculated from the percentile scores of the raw scores on the common

items for each measure. A set of  $z$ -scores on the common items of each measure is selected for linking the two measures (ideally 10–20  $z$ -scores on the common items of each measure that are not in the extremes of the distribution to prevent a distorted transformed scale). The same number of  $z$ -scores are selected from the common items of each measure to generate  $z$ -score pairs for linking the two measures (e.g., in the present study, we selected 17  $z$ -scores from the common items of each measure resulting in 17  $z$ -score pairs). The first assumption of Thurstone scaling in a common-item design is that the association between these  $z$ -score pairs (i.e., the selected set of  $z$ -scores for the common items) of the two measures to be linked is linear (Kolen & Brennan, 2014). The second assumption of Thurstone scaling is that the underlying trait is normally distributed. After examining the assumptions of Thurstone scaling, we linked YASR scores at age 20 to YSR scores at age 19 (i.e., the target scale that serves as the anchor). To do this, we applied the following steps, separately for males and females<sup>1</sup>, of Thurstone scaling in a common-item design (Kolen & Brennan, 2014) to link YASR scores with YSR scores:

(1) To ensure a meaningful mean-level of change of internalizing problem scores across ages 14–24, we first examined scores on the 17 common items (i.e., the items that were common to both the YSR and YASR). The raw frequency distribution of scores on the common items is in Tables S1–S2). Participants' mean scores on the common items are depicted in Panel A of Figure S2. The percentile ranks of raw scores on the common items are in Tables S3–S4.

(2) For both ages, we calculated  $z$ -scores of raw scores on the common items within age based on the percentile ranks from step 1, see Tables S5–S6.

---

<sup>1</sup> This was done because of the robust sex differences in levels of internalizing problems among adolescents and young adults, with females having higher levels than males (Hankin et al., 1998).

(3) For both ages, we calculated the mean and standard deviation of the unique  $z$ -scores (i.e., the  $z$ -score values in a given column of Tables S5–S6)<sup>3</sup> based on those unique raw scores whose associated  $z$ -scores were between  $-2$  and  $+2$  at the age of the target scale (i.e., age 19). We selected  $z$ -scores between  $-2$  and  $+2$  as recommended by Kolen and Brennan (2014) and because trimmed  $z$ -scores are more accurate in Thurstone scaling than using all  $z$ -scores. We calculated the *population* standard deviation of the  $z$ -scores, consistent with Kolen and Brennan (2014); all descriptive statistics of the sample used the *sample* standard deviation. The calculated mean and standard deviation of the  $z$ -scores between  $-2$  and  $+2$  at the age of the target scale are at the bottom of Tables S5–S6.

(4) The mean and standard deviation of the scaled scores were calculated by the following formulas (adapted from Kolen & Brennan, 2014):

$$\mu(\text{YASR}) = \sigma(\text{YSR}) \left[ \mu(z_{\text{YSR}}) - \frac{\sigma(z_{\text{YSR}})}{\sigma(z_{\text{YASR}})} \mu(z_{\text{YASR}}) \right] + \mu(\text{YSR}) \quad (\text{S1})$$

$$\sigma(\text{YASR}) = \frac{\sigma(z_{\text{YSR}})}{\sigma(z_{\text{YASR}})} \sigma(\text{YSR}) \quad (\text{S2})$$

where YASR and YSR are vectors of raw scores on the YASR or YSR, respectively;  $z_{\text{YASR}}$  and  $z_{\text{YSR}}$  are the unique  $z$ -scores based on those raw scores whose associated  $z$ -scores were between  $-2$  and  $+2$  at age 19 (from step 3). The first component of Formula S1 scales the mean of the  $z$ -scores of the common items at age 20 to be on a  $z$ -score metric relative to the  $z$ -score metric of the common items at age 19 in order to retain changes in means and variances from age 19 to 20:  $\left[ \mu(z_{\text{YSR}}) - \frac{\sigma(z_{\text{YSR}})}{\sigma(z_{\text{YASR}})} \mu(z_{\text{YASR}}) \right]$ . The second component of Formula S1 then multiplies the scaled  $z$ -score metric by the standard deviation of the target scale and adds the mean of the target scale. This re-scales the age 20  $z$ -score metric (on a scale that is relative to the age 19  $z$ -score metric) to the metric of the total raw score at age 19 in order to make the re-scaled scores at

age 20 comparable to the raw scores at age 19. Thus, the YASR scores at age 20 are re-scaled to be on the scale of the YSR at age 19, while still retaining changes in means and variances over time (based on the changes in means and variances of the common items).

Consistent with recommendations (Kolen & Brennan, 2014), we then applied a process of linking and chaining to link the remaining YASR scores to the YSR metric at age 19 based on the raw frequency distribution (Tables S7–S8), percentile ranks (Tables S9–S10), and  $z$ -scores of the total raw scores (Tables S11–S12). To do so, we repeated steps 1–4 by (a) linking the YASR scores at age 21 to the newly scaled YASR scores at age 20, (b) linking scores at age 22 to the newly scaled scores at age 21, (c) linking scores at age 23 to the newly scaled scores at age 22, and (d) linking scores at age 24 to the newly scaled scores at age 23. Linking and chaining allowed us to calculate a mean and standard deviation for the scaled YASR score at each age from ages 20 to 24. Based on the mean and standard deviation of the scaled scores for each year, we calculated a conversion table for converting YASR scores to YSR equivalents based on the scale of the YSR scores at age 19. We calculated a conversion table by multiplying the  $z$ -scores of the total raw scores (Tables S11–S12) by the standard deviation of the scaled score (from Equation S2; bottom of Tables S5–S6) and added the mean of the scaled score (from Equation S1; bottom of Tables S5–S6). The conversion table for converting YASR scores to YSR equivalents is in Table S13.

## **Results**

**Assumptions of Thurstone Scaling.** We examined the two assumptions of Thurstone scaling in a common-item design: (1) the association between the selected  $z$ -score pairs from the common items of the two measures to be linked is linear, and (2) the underlying trait is normally distributed. Regarding assumption 1, we observed that the associations between the selected  $z$ -



score pairs from the common items of adjacent years were highly linear; no curvilinearity was observed. Regarding assumption 2, we observed that the raw total Internalizing scores were positively skewed (skew values ranged from 0.93 to 1.72 across years). Despite the skewed scores, it is plausible that the underlying trait (i.e., the internalizing spectrum) is normally distributed, especially given evidence that internalizing problems are dimensionally rather than categorically distributed (Markon, Chmielewski, & Miller, 2011). That is, the latent trait of internalizing is likely normally distributed even if the observed scores are not, which would be consistent with the assumption of Thurstone scaling. Thus, given theoretical and empirical evidence supporting that we approximately met these assumptions, we proceeded with vertical scaling using Thurstone scaling.

**Linking the YASR Scores to the Scale of the YSR at Age 19.** Next, we linked the YASR and YSR scores so that scores on the two measures were on the same scale and could be compared. To link the two measures, we re-scaled the YASR scores at age 20 to the scale of the YSR at age 19 (see steps 1–4 from the Statistical Analysis section of the Method section). The Thurstone scaling approach to vertical scaling is depicted in Figure S1. First, we examined scores at ages 19 and 20 on the 17 common items (i.e., the items that were common to both the YSR and YASR), and calculated percentile ranks, see Tables S3–S4. Second, we calculated  $z$ -scores of raw scores on the common items within age (Tables S5–S6) based on the percentile ranks from step 1. Third, we calculated the mean and standard deviation of the unique  $z$ -scores based on those raw scores (males: 0–15; females: 0–16) whose associated  $z$ -scores were between -2 and +2 at the age of the target scale (i.e., age 19).<sup>2</sup> The mean and standard deviation of these

---

<sup>2</sup> That is, we used the *unique*  $z$ -scores, not the vector of all  $z$ -scores from the raw scores (multiple participants may have the same raw score and therefore the same  $z$ -score).

$z$ -scores are at the bottom of Tables S5–S6. Fourth, we calculated the mean and standard deviation of the scaled scores on the scale of the YSR at age 19 based on the mean and standard deviation of the raw scores and selected  $z$ -scores at each age.

The mean of the scaled score at age 20 for females was calculated as (equation S1):

$$\mu(\text{YASR}_{20}) = 7.40 \left[ 0.24 - \frac{0.96}{0.93} 0.26 \right] + 9.70 = 9.50.^3$$

The standard deviation of the scaled score at age 20 for females was calculated as (equation S2):  $\sigma(\text{YASR}_{20}) = \frac{0.96}{0.93} 7.40 = 7.63.^4$

We then applied linking and chaining to link the remaining YASR scores to the YSR metric at age 19. To do so, we repeated steps 1–4 by linking the YASR scores at age 21 to the newly scaled YASR scores at age 20, linking scores at age 22 to the newly scaled scores at age 21, etc.

**Conversion Table for Converting YASR Scores to YSR Equivalents.** Linking and chaining allowed us to calculate a mean and standard deviation for the scaled YASR score at each age from ages 20 to 24. Based on the mean and standard deviation of the scaled scores for each year, we converted YASR scores to YSR equivalents based on the scale of the YSR scores at age 19. To do this, we multiplied the  $z$ -scores of the total raw scores (Tables S11–S12) by the standard deviation of the scaled score (from Equation S2; bottom of Table S13) and added the mean of the scaled score (from Equation S1; bottom of Table S13).

The conversion table for converting YASR scores to YSR equivalents in our sample is in Table S13. Visual examination of the conversion table shows that many of the scores were

---

<sup>3</sup> These values came from the following sources: 7.40 (Table S7), 0.24 (Table S5), 0.96 (Table S5), 0.93 (Table S5), 0.26 (Table S5), 9.70 (Table S7), 9.50 (Table 2). Note that the above calculations are slightly different from the actual calculations due to rounding error.

<sup>4</sup> These values came from the following sources: 0.96 (Table S5), 0.93 (Table S5), 7.40 (Table S7), 7.63 (Table 2). Note that the above calculations are slightly different from the actual calculations due to rounding error.

highly similar before and after rescaling, while rescaling changed some of the scores by more than 2 points (particularly for females). The mean and standard deviation of the scaled scores are at the bottom of Table 13. Participants' mean internalizing problem scores, after rescaling the YASR scores to the metric of the YSR, are depicted in Panel B of Figure S2. Notably, the scores retained a highly similar pattern of mean-level change when examining the re-scaled total scores compared to when examining just the common items (see Panel A of Figure S2). Thus, the Thurstone Scaling approach successfully retained mean-level change when re-scaling the YASR scores to be on the same metric as the YSR while still using a more comparable scale.

**Growth Curve Model.** To examine growth curves, we first compared a linear growth curve model to a quadratic growth curve model in HLM to identify the best-fitting form of change for the rescaled internalizing problem scores. The model that allowed quadratic slopes to vary across individuals (i.e., random quadratic slopes) was not positive definite (i.e., not all variances in the variance-covariance matrix were non-zero and positive), likely because the variance in the quadratic term was close to zero ( $\tau_{22} < .0001$ ). The small variance in the quadratic term suggested that individuals did not significantly differ in quadratic curvature. A model with random linear slopes and a quadratic term that was fixed across individuals (i.e., fixed quadratic slopes) fit better than a model with only random linear slopes ( $\chi^2[1] = 24.46, p < .001$ ). A model with fixed cubic slopes fit better than the model with random linear slopes and fixed quadratic slopes ( $\chi^2[1] = 5.63, p = .018$ ). A model with fixed quartic slopes fit better than the cubic model ( $\chi^2[1] = 4.65, p = .031$ ), and was the best-fitting model (a model with fixed fifth-degree polynomial slopes did not significantly improve fit;  $\chi^2[1] = 3.07, p = .080$ ). Individuals' quartic trajectories, and the average quartic trajectory for males and females are depicted in Figure S3. The average quartic trajectory showed slight decreases over time,

primarily for females.

Overall, the growth curves showed little curvature, which would be consistent with evidence that likelihood ratio tests may be sensitive to small fit differences with larger sample sizes (Tomarken & Waller, 2003). Thus, the polynomial growth terms may have over-fit the data, especially given the lengthy developmental span. Moreover, there are difficulties in interpreting and replicating findings from polynomial growth models, and mapping polynomial growth terms onto developmental theory (Grimm, Ram, & Hamagami, 2011). For these reasons, for comparing the common items to the rescaled scores and for examining the predictors of change in internalizing problems, we examined the general form of change by examining the linear model for ease of interpretation.

In the linear growth curve model with no predictors of the intercepts or slopes, intercepts reflected an individual's estimated initial level of internalizing problems at age 14. Slopes reflected participants' linear change in internalizing problems over time. There was evidence of a negative slope ( $B = -0.08$ ,  $t[3980] = -1.94$ ,  $p = .053$ ). In a similar growth curve model examining the trajectories of scores on the common items, however, the slope was not significant ( $B = -0.03$ ,  $t[3980] = -1.08$ ,  $p = .282$ ).

Although the form of change for the age-relevant items versus common items was highly similar at the *group-level*, there were differences at the *individual-level*. Some participants showed *decreases* in internalizing problems over time when using the age-relevant items while they showed *increases* in internalizing problems when using the common items (or vice versa). The participants who showed decreases using the age-relevant items and increases using the common items presumably had higher levels of internalizing problems on the *non-common* items of the YSR (i.e., items that were on the Internalizing scale of the YSR but not the YASR) or

lower levels on the *non-common* items of the YASR (compared to the other participants).

Because the Somatic Complaints subscale was included in the Internalizing Scale of the YSR but not YASR, the majority (9 items, 60%) of the non-common Internalizing items of the YSR were items assessing somatic complaints. Therefore, we examined participants' levels of somatic complaints on the YSR. Consistent with expectations, participants who showed decreases in internalizing problems using the age-relevant items but increases using the common items showed higher mean levels of somatic complaints from ages 14–19 ( $M = 3.22$ ) than participants who did not ( $M = 1.81$ ;  $t[32.38] = -3.41$ ,  $p = .002$ ). The reverse was also true; participants who showed increases in internalizing problems using the age-relevant items but decreases using the common items, showed lower mean levels of somatic complaints from ages 14–19 ( $M = 0.93$ ) than participants who did not ( $M = 1.94$ ;  $t[30.76] = 4.49$ ,  $p < .001$ ).

We then examined sex and ethnicity as predictors of the intercepts and linear slopes of the rescaled internalizing problem scores (see Table S14). There were no significant linear slopes when controlling for the other model predictors. Females showed higher intercepts than males, but males and females did not significantly differ in their linear slopes. African Americans and those of “other” ethnicity did not significantly differ from European Americans in their intercepts or linear slopes. The model accounted for approximately three-fourths of the variance in internalizing problems over time.



Table S1. Raw score frequency distributions on common items for females.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	6	10	8	5	12	17	14	15	20	21
1	12	8	8	18	13	18	25	13	23	18
2	19	19	21	18	30	27	20	17	12	18
3	20	18	23	23	20	21	15	19	19	19
4	18	17	15	16	20	20	23	21	17	14
5	13	12	23	14	12	14	15	15	22	21
6	18	14	15	20	11	19	27	25	15	19
7	16	15	16	17	28	18	12	17	20	7
8	13	12	11	8	10	18	11	15	12	13
9	11	14	10	10	15	14	18	18	18	19
10	11	9	10	6	20	9	14	9	14	20
11	9	13	15	7	6	10	10	8	9	6
12	11	4	7	10	5	7	6	9	8	6
13	6	6	9	11	1	5	8	7	5	11
14	7	9	5	5	6	4	2	2	10	12
15	7	2	7	6	3	6	5	6	5	5
16	5	6	2	2	1	2	2	2	6	3
17	3	7	4	7	4	3	3	5	3	2
18	5	2	5	2	5	2	1	4	2	2
19	1	5	4	2	1	4	3	3	4	5
20	1	1	2	1	4	3	1	2	1	1
21	0	2	1	2	1	0	1	0	0	2
22	0	4	0	2	1	0	2	2	0	1
23	0	0	3	1	1	0	2	1	1	1
24	0	1	1	1	0	1	0	0	3	1
25	0	1	0	0	0	2	0	0	0	2
26	0	0	0	0	0	0	1	0	0	0
27	0	0	0	0	0	0	1	1	0	0
28	1	0	1	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0
<i>M</i>	7.42	8.05	7.80	7.42	6.77	6.64	6.79	7.16	7.06	7.34
<i>SD</i>	5.03	5.77	5.58	5.45	5.09	5.25	5.37	5.25	5.35	5.62

Note: mean and standard deviation reflect the mean and standard deviation of the participants' scores on the common items (they do not reflect the mean and standard deviation of the values in the above column).

Table S2. Raw score frequency distributions on common items for males.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	16	16	17	22	18	27	33	27	34	26
1	27	24	24	23	26	30	18	30	22	20
2	24	21	26	33	29	22	25	25	31	21
3	24	21	21	22	27	23	26	24	18	23
4	17	11	25	23	26	21	18	18	19	16
5	11	23	27	17	18	20	15	11	18	13
6	10	18	14	13	12	14	11	14	15	11
7	17	11	10	10	17	10	15	12	12	19
8	13	6	11	5	12	9	11	12	8	10
9	6	4	8	11	8	10	7	8	9	10
10	5	6	9	6	8	6	6	10	6	4
11	6	7	8	9	7	8	6	5	10	6
12	7	8	4	3	2	8	7	7	6	4
13	3	6	9	2	6	9	4	8	10	10
14	4	4	4	3	3	4	3	2	5	3
15	5	4	1	1	5	2	1	6	2	1
16	1	0	1	1	1	4	4	1	1	3
17	0	1	0	2	2	2	6	3	1	3
18	1	1	3	1	3	3	3	1	5	3
19	0	1	2	5	1	1	4	4	3	2
20	0	1	1	0	2	0	0	1	1	2
21	0	0	0	1	0	0	0	0	0	2
22	1	1	0	1	0	0	0	0	0	1
23	1	1	1	1	0	0	0	0	0	1
24	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	1	0	0	0	0
26	0	0	0	0	0	1	0	0	0	1
27	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	1	0	0	0	0	0
29	0	0	0	0	0	0	0	1	1	0
<i>M</i>	5.31	5.63	5.54	5.24	5.49	5.58	5.52	5.59	5.59	6.11
<i>SD</i>	4.52	4.76	4.54	4.89	4.76	5.03	5.04	5.15	5.17	5.53

Note: mean and standard deviation reflect the mean and standard deviation of the participants' scores on the common items (they do not reflect the mean and standard deviation of the values in the above column).



Table S4. Percentile ranks (divided by 100) on common items for males.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	.04	.04	.04	.05	.04	.06	.08	.06	.07	.06
1	.15	.14	.13	.16	.13	.18	.19	.18	.19	.17
2	.28	.26	.24	.29	.25	.29	.29	.30	.30	.27
3	.40	.37	.35	.41	.37	.39	.40	.41	.41	.37
4	.50	.45	.45	.52	.48	.48	.50	.50	.49	.46
5	.57	.54	.56	.61	.58	.57	.57	.57	.56	.53
6	.62	.64	.65	.68	.64	.64	.63	.62	.63	.58
7	.69	.71	.70	.73	.71	.69	.69	.67	.69	.65
8	.77	.76	.75	.77	.76	.73	.75	.73	.73	.72
9	.81	.78	.79	.81	.81	.77	.79	.77	.77	.76
10	.84	.81	.83	.85	.84	.80	.82	.81	.80	.80
11	.87	.84	.87	.88	.88	.83	.84	.84	.83	.82
12	.90	.88	.89	.91	.89	.87	.87	.87	.86	.84
13	.93	.91	.92	.92	.91	.91	.90	.90	.90	.87
14	.94	.94	.95	.93	.93	.93	.91	.92	.93	.91
15	.97	.96	.96	.94	.95	.94	.92	.94	.95	.92
16	.98	–	.97	.95	.96	.96	.93	.96	.95	.93
17	–	.97	–	.95	.97	.97	.96	.97	.96	.94
18	.99	.98	.98	.96	.98	.98	.98	.97	.97	.95
19	–	.98	.99	.98	.99	.99	1.00	.98	.99	.96
20	–	.99	.99	–	.99	–	–	.99	.99	.97
21	–	–	–	.99	–	–	–	–	–	.98
22	.99	.99	–	.99	–	–	–	–	–	.99
23	1.00	1.00	1.00	1.00	–	–	–	–	–	.99
24	–	–	–	–	–	–	–	–	–	–
25	–	–	–	–	–	.99	–	–	–	–
26	–	–	–	–	–	1.00	–	–	–	1.00
27	–	–	–	–	–	–	–	–	–	–
28	–	–	–	–	1.00	–	–	–	–	–
29	–	–	–	–	–	–	–	1.00	1.00	–

Table S5. Z-scores of common items for females.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	-1.47	-1.40	-1.40	-1.36	-1.33	-1.26	-1.27	-1.36	-1.32	-1.31
1	-1.28	-1.22	-1.22	-1.18	-1.13	-1.07	-1.08	-1.17	-1.13	-1.13
2	-1.08	-1.05	-1.04	-0.99	-0.94	-0.88	-0.89	-0.98	-0.95	-0.95
3	-0.88	-0.88	-0.86	-0.81	-0.74	-0.69	-0.71	-0.79	-0.76	-0.77
4	-0.68	-0.70	-0.68	-0.63	-0.54	-0.50	-0.52	-0.60	-0.57	-0.59
5	-0.48	-0.53	-0.50	-0.44	-0.35	-0.31	-0.33	-0.41	-0.39	-0.42
6	-0.28	-0.36	-0.32	-0.26	-0.15	-0.12	-0.15	-0.22	-0.20	-0.24
7	-0.08	-0.18	-0.14	-0.08	0.05	0.07	0.04	-0.03	-0.01	-0.06
8	0.12	-0.01	0.04	0.11	0.24	0.26	0.22	0.16	0.17	0.12
9	0.31	0.16	0.21	0.29	0.44	0.45	0.41	0.35	0.36	0.30
10	0.51	0.34	0.39	0.47	0.63	0.64	0.60	0.54	0.55	0.47
11	0.71	0.51	0.57	0.66	0.83	0.83	0.78	0.73	0.74	0.65
12	0.91	0.68	0.75	0.84	1.03	1.02	0.97	0.92	0.92	0.83
13	1.11	0.86	0.93	1.02	1.22	1.21	1.16	1.11	1.11	1.01
14	1.31	1.03	1.11	1.21	1.42	1.40	1.34	1.30	1.30	1.18
15	1.51	1.20	1.29	1.39	1.62	1.59	1.53	1.49	1.48	1.36
16	1.71	1.38	1.47	1.57	1.81	1.78	1.72	1.68	1.67	1.54
17	1.91	1.55	1.65	1.76	2.01	1.98	1.90	1.87	1.86	1.72
18	2.10	1.72	1.83	1.94	2.21	2.17	2.09	2.07	2.04	1.90
19	2.30	1.90	2.01	2.12	2.40	2.36	2.28	2.26	2.23	2.07
20	2.50	2.07	2.18	2.31	2.60	2.55	2.46	2.45	2.42	2.25
21	–	2.24	2.36	2.49	2.80	–	2.65	–	–	2.43
22	–	2.42	–	2.67	2.99	–	2.83	2.83	–	2.61
23	–	–	2.72	2.86	3.19	–	3.02	3.02	2.98	2.79
24	–	2.76	2.90	3.04	–	3.31	–	–	3.16	2.96
25	–	2.94	–	–	–	3.50	–	–	–	3.14
26	–	–	–	–	–	–	3.58	–	–	–
27	–	–	–	–	–	–	3.77	3.78	–	–
28	4.09	–	3.62	–	–	–	–	–	–	–
29	–	–	–	–	–	–	–	–	–	–
<i>M</i>	0.12	-0.01	0.04	0.11	0.24	0.26	0.22	0.16	0.17	0.12
<i>SD</i>	0.97	0.85	0.88	0.90	0.96	0.93	0.91	0.93	0.92	0.87

Note: the dashed line reflects those raw scores at age 19 whose associated  $z$ -scores were between -2 and +2 (i.e., raw scores of 0 to 16). Mean and standard deviation reflect the mean and standard deviation of the  $z$ -scores whose associated raw scores ranged from 0 to 16 (i.e., the mean and standard deviation of the values above the dashed line).

Table S6. Z-scores of common items for males.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	-1.18	-1.18	-1.22	-1.07	-1.15	-1.11	-1.10	-1.09	-1.08	-1.11
1	-0.95	-0.97	-1.00	-0.87	-0.94	-0.91	-0.90	-0.89	-0.89	-0.92
2	-0.73	-0.76	-0.78	-0.66	-0.73	-0.71	-0.70	-0.70	-0.69	-0.74
3	-0.51	-0.55	-0.56	-0.46	-0.52	-0.51	-0.50	-0.50	-0.50	-0.56
4	-0.29	-0.34	-0.34	-0.25	-0.31	-0.31	-0.30	-0.31	-0.31	-0.38
5	-0.07	-0.13	-0.12	-0.05	-0.10	-0.12	-0.10	-0.11	-0.11	-0.20
6	0.15	0.08	0.10	0.16	0.11	0.08	0.10	0.08	0.08	-0.02
7	0.38	0.29	0.32	0.36	0.32	0.28	0.29	0.27	0.27	0.16
8	0.60	0.50	0.54	0.57	0.53	0.48	0.49	0.47	0.47	0.34
9	0.82	0.71	0.76	0.77	0.74	0.68	0.69	0.66	0.66	0.52
10	1.04	0.92	0.98	0.97	0.95	0.88	0.89	0.86	0.85	0.70
11	1.26	1.13	1.20	1.18	1.16	1.08	1.09	1.05	1.05	0.88
12	1.48	1.34	1.42	1.38	1.37	1.28	1.29	1.25	1.24	1.06
13	1.70	1.55	1.64	1.59	1.58	1.47	1.49	1.44	1.43	1.25
14	1.93	1.76	1.86	1.79	1.79	1.67	1.68	1.63	1.63	1.43
15	2.15	1.97	2.08	2.00	2.00	1.87	1.88	1.83	1.82	1.61
16	2.37	–	2.30	2.20	2.21	2.07	2.08	2.02	2.01	1.79
17	–	2.39	–	2.41	2.42	2.27	2.28	2.22	2.21	1.97
18	2.81	2.60	2.74	2.61	2.63	2.47	2.48	2.41	2.40	2.15
19	–	2.81	2.96	2.82	2.84	2.67	2.68	2.61	2.59	2.33
20	–	3.02	3.18	–	3.05	–	–	2.80	2.79	2.51
21	–	–	–	3.23	–	–	–	–	–	2.69
22	3.70	3.44	–	3.43	–	–	–	–	–	2.87
23	3.92	3.65	3.84	3.64	–	–	–	–	–	3.05
24	–	–	–	–	–	–	–	–	–	–
25	–	–	–	–	–	3.86	–	–	–	–
26	–	–	–	–	–	4.06	–	–	–	3.60
27	–	–	–	–	–	–	–	–	–	–
28	–	–	–	–	4.72	–	–	–	–	–
29	–	–	–	–	–	–	–	4.55	4.53	–
<i>M</i>	0.49	0.39	0.43	0.46	0.42	0.38	0.39	0.37	0.37	0.25
<i>SD</i>	1.02	0.97	1.01	0.94	0.97	0.92	0.92	0.90	0.89	0.83

Note: the dashed line reflects those raw scores at age 19 whose associated  $z$ -scores were between -2 and +2 (i.e., raw scores of 0 to 15). Mean and standard deviation reflect the mean and standard deviation of the  $z$ -scores whose associated raw scores ranged from 0 to 15 (i.e., the mean and standard deviation of the values above the dashed line).

Table S7. Raw score frequency distributions on all items for females.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	5	9	5	4	8	6	7	8	12	11
1	10	4	2	10	8	19	15	9	14	13
2	11	9	12	14	19	13	12	13	14	16
3	14	14	12	15	16	21	19	15	15	17
4	12	10	18	17	21	16	20	16	8	13
5	4	10	11	8	11	9	12	10	13	13
6	15	20	12	12	11	15	10	17	21	14
7	13	11	13	17	9	20	20	13	15	15
8	16	8	12	10	12	13	16	17	7	9
9	8	7	14	11	14	16	11	13	12	10
10	14	10	13	7	12	9	11	13	11	8
11	8	6	7	11	11	13	8	12	18	9
12	14	9	16	9	13	12	13	11	11	11
13	7	6	13	5	12	15	14	7	14	13
14	4	10	7	8	5	6	7	15	8	16
15	8	6	5	6	8	1	4	7	6	9
16	5	7	3	4	3	6	6	3	5	8
17	6	6	6	7	3	4	6	5	6	4
18	2	2	5	5	4	5	5	5	8	7
19	3	3	4	0	4	3	3	2	6	7
20	5	10	1	4	4	4	5	4	4	4
21	4	4	4	3	3	3	2	5	2	4
22	9	4	4	7	2	3	3	2	4	1
23	2	3	4	2	2	1	2	1	4	2
24	0	3	8	4	0	1	0	4	1	2
25	3	4	1	1	2	3	0	2	3	2
26	1	3	2	1	4	2	4	3	2	3
27	3	1	4	0	2	1	1	0	1	2
28	0	1	1	3	2	1	0	1	0	2
29	4	2	1	2	1	0	1	1	0	0
30	0	3	0	1	0	0	2	0	0	1
31	0	1	0	1	1	0	0	0	3	0
32	0	2	0	3	1	1	0	0	0	1
33	2	1	0	0	1	1	1	0	0	0
34	0	2	1	1	1	1	0	1	1	1
35	0	0	0	1	0	0	1	0	0	1
36	0	0	2	0	0	0	1	0	0	0
37	0	0	1	0	0	0	0	1	0	0
38	0	0	1	0	0	0	0	0	0	0
39	0	0	0	0	0	0	0	0	0	0
40	0	0	0	0	0	0	0	0	0	0
41	0	0	1	0	0	0	0	0	0	0
42	0	0	0	0	0	0	0	0	0	0
43	0	0	0	0	0	0	0	0	0	0
44	0	0	0	0	0	0	0	0	0	0
45	0	0	0	0	0	0	0	0	0	0
46	0	0	0	0	0	0	0	0	0	0
47	1	0	0	0	0	0	0	0	0	0
<i>M</i>	10.80	11.64	11.38	10.55	9.70	9.07	9.42	9.77	9.84	10.07
<i>SD</i>	7.85	8.38	8.09	8.03	7.40	6.83	7.13	6.90	7.10	7.42

Note: mean and standard deviation reflect the mean and standard deviation of the participants' scores on all items (they do not reflect the mean and standard deviation of the values in the above column).



Table S8. Raw score frequency distributions on all items for males.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	9	11	13	17	13	13	14	9	15	15
1	14	16	18	18	17	16	18	18	24	10
2	20	18	15	24	18	24	20	29	17	17
3	21	14	19	19	26	28	20	20	14	18
4	15	14	19	19	18	8	16	16	29	16
5	5	21	16	15	19	14	14	16	14	19
6	12	11	16	10	16	24	14	11	10	12
7	17	9	13	14	13	14	15	9	14	11
8	9	11	16	11	15	7	9	13	9	10
9	12	6	8	7	11	13	11	14	12	14
10	9	9	7	6	10	10	8	6	11	7
11	11	8	10	5	4	6	8	11	6	11
12	6	7	9	7	7	5	6	5	5	3
13	7	3	9	4	11	5	9	5	6	3
14	4	6	4	3	5	5	4	5	7	7
15	4	3	3	7	4	5	4	6	6	3
16	2	3	5	1	1	7	6	5	4	3
17	2	5	4	4	4	10	3	2	6	4
18	4	4	7	6	2	3	1	6	8	5
19	0	2	2	3	1	3	1	10	4	6
20	3	2	4	1	5	2	5	1	2	3
21	3	1	0	0	2	2	1	5	1	0
22	1	3	1	2	3	1	4	0	2	2
23	4	1	1	2	2	2	4	0	2	4
24	0	1	1	2	3	2	3	2	1	1
25	0	2	1	1	0	2	2	2	5	1
26	0	1	2	1	1	1	2	1	0	2
27	0	0	0	1	0	1	0	0	0	2
28	2	1	1	3	0	0	1	2	2	1
29	0	0	0	0	1	0	0	0	0	2
30	2	1	1	1	1	0	0	0	0	1
31	0	1	0	0	0	0	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0
33	0	0	0	0	0	0	0	0	0	0
34	1	1	1	0	0	0	0	0	0	1
35	0	0	0	0	1	1	0	0	0	0
36	0	0	0	0	0	0	0	0	0	0
37	0	0	0	1	0	1	0	0	0	1
38	0	0	0	0	0	0	0	0	1	0
39	0	0	0	0	0	0	0	0	0	0
40	0	0	0	0	0	0	0	1	0	0
<i>M</i>	7.98	8.02	7.89	7.52	7.65	7.94	7.95	8.10	8.04	8.74
<i>SD</i>	6.66	6.93	6.45	7.12	6.51	6.80	6.73	6.86	6.94	7.55

Note: mean and standard deviation reflect the mean and standard deviation of the participants' scores on all items (they do not reflect the mean and standard deviation of the values in the above column).



Table S10. Percentile ranks (divided by 100) on all items for males.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	.03	.03	.03	.04	.03	.03	.03	.02	.03	.04
1	.08	.10	.10	.12	.09	.09	.10	.08	.11	.09
2	.17	.18	.17	.22	.17	.17	.19	.18	.20	.16
3	.27	.27	.25	.32	.26	.29	.28	.29	.27	.24
4	.36	.34	.33	.41	.35	.36	.36	.37	.36	.32
5	.41	.43	.41	.49	.44	.41	.43	.43	.45	.40
6	.45	.51	.48	.54	.51	.49	.49	.50	.50	.47
7	.53	.56	.54	.60	.57	.57	.56	.54	.55	.53
8	.59	.61	.61	.66	.63	.62	.61	.59	.60	.57
9	.64	.65	.66	.70	.69	.66	.65	.64	.64	.63
10	.70	.69	.69	.73	.73	.71	.70	.69	.69	.68
11	.75	.73	.73	.76	.76	.74	.73	.73	.73	.72
12	.79	.78	.77	.79	.79	.77	.76	.76	.75	.75
13	.82	.80	.81	.81	.82	.79	.80	.78	.77	.77
14	.85	.82	.84	.83	.86	.81	.83	.80	.80	.79
15	.87	.85	.86	.85	.88	.83	.84	.83	.83	.81
16	.88	.86	.88	.87	.89	.86	.87	.85	.85	.83
17	.89	.88	.89	.88	.90	.89	.89	.87	.87	.84
18	.91	.90	.92	.90	.91	.92	.90	.88	.90	.87
19	–	.92	.94	.93	.92	.93	.90	.92	.92	.89
20	.93	.93	.95	.93	.93	.94	.91	.94	.94	.91
21	.94	.94	–	–	.94	.95	.93	.96	.95	–
22	.95	.95	.96	.94	.96	.96	.94	–	.95	.92
23	.96	.96	.97	.95	.97	.96	.96	–	.96	.93
24	–	.96	.97	.96	.98	.97	.97	.97	.97	.95
25	–	.97	.98	.97	–	.98	.98	.98	.98	.95
26	–	.98	.98	.97	.99	.99	.99	.99	–	.96
27	–	–	–	.98	–	.99	–	–	–	.97
28	.98	.98	.99	.99	–	–	1.00	.99	.99	.98
29	–	–	–	–	.99	–	–	–	–	.98
30	.99	.99	.99	.99	.99	–	–	–	–	.99
31	–	.99	–	–	–	–	–	–	–	–
32	–	–	–	–	–	–	–	–	–	–
33	–	–	–	–	–	–	–	–	–	–
34	1.00	1.00	1.00	–	–	–	–	–	–	.99
35	–	–	–	–	1.00	.99	–	–	–	–
36	–	–	–	–	–	–	–	–	–	–
37	–	–	–	1.00	–	1.00	–	–	–	1.00
38	–	–	–	–	–	–	–	–	1.00	–
39	–	–	–	–	–	–	–	–	–	–
40	–	–	–	–	–	–	–	1.00	–	–



Table S12. Z-scores of all items for males.

score	Age (years)									
	14	15	16	17	19	20	21	22	23	24
0	-1.20	-1.16	-1.22	-1.06	-1.18	-1.17	-1.18	-1.18	-1.16	-1.16
1	-1.05	-1.01	-1.07	-0.92	-1.02	-1.02	-1.03	-1.03	-1.02	-1.02
2	-0.90	-0.87	-0.91	-0.78	-0.87	-0.87	-0.88	-0.89	-0.87	-0.89
3	-0.75	-0.72	-0.76	-0.64	-0.71	-0.73	-0.74	-0.74	-0.73	-0.76
4	-0.60	-0.58	-0.60	-0.49	-0.56	-0.58	-0.59	-0.60	-0.58	-0.63
5	-0.45	-0.44	-0.45	-0.35	-0.41	-0.43	-0.44	-0.45	-0.44	-0.50
6	-0.30	-0.29	-0.29	-0.21	-0.25	-0.29	-0.29	-0.31	-0.29	-0.36
7	-0.15	-0.15	-0.14	-0.07	-0.10	-0.14	-0.14	-0.16	-0.15	-0.23
8	0.00	0.00	0.02	0.07	0.05	0.01	0.01	-0.01	-0.01	-0.10
9	0.15	0.14	0.17	0.21	0.21	0.16	0.16	0.13	0.14	0.03
10	0.30	0.29	0.33	0.35	0.36	0.30	0.31	0.28	0.28	0.17
11	0.45	0.43	0.48	0.49	0.51	0.45	0.45	0.42	0.43	0.30
12	0.60	0.57	0.64	0.63	0.67	0.60	0.60	0.57	0.57	0.43
13	0.75	0.72	0.79	0.77	0.82	0.74	0.75	0.71	0.71	0.56
14	0.90	0.86	0.95	0.91	0.98	0.89	0.90	0.86	0.86	0.70
15	1.05	1.01	1.10	1.05	1.13	1.04	1.05	1.01	1.00	0.83
16	1.20	1.15	1.26	1.19	1.28	1.18	1.20	1.15	1.15	0.96
17	1.35	1.30	1.41	1.33	1.44	1.33	1.35	1.30	1.29	1.09
18	1.50	1.44	1.57	1.47	1.59	1.48	1.49	1.44	1.44	1.23
19	–	1.58	1.72	1.61	1.74	1.63	1.64	1.59	1.58	1.36
20	1.80	1.73	1.88	1.75	1.90	1.77	1.79	1.73	1.72	1.49
21	1.95	1.87	–	–	2.05	1.92	1.94	1.88	1.87	–
22	2.10	2.02	2.19	2.03	2.20	2.07	2.09	–	2.01	1.76
23	2.26	2.16	2.34	2.17	2.36	2.21	2.24	–	2.16	1.89
24	–	2.30	2.50	2.32	2.51	2.36	2.39	2.32	2.30	2.02
25	–	2.45	2.65	2.46	–	2.51	2.54	2.46	2.44	2.15
26	–	2.59	2.81	2.60	2.82	2.66	2.68	2.61	–	2.29
27	–	–	–	2.74	–	2.80	–	–	–	2.42
28	3.01	2.88	3.12	2.88	–	–	2.98	2.90	2.88	2.55
29	–	–	–	–	3.28	–	–	–	–	2.68
30	3.31	3.17	3.42	3.16	3.43	–	–	–	–	2.81
31	–	3.31	–	–	–	–	–	–	–	–
32	–	–	–	–	–	–	–	–	–	–
33	–	–	–	–	–	–	–	–	–	–
34	3.91	3.75	4.04	–	–	–	–	–	–	3.34
35	–	–	–	–	4.20	3.98	–	–	–	–
36	–	–	–	–	–	–	–	–	–	–
37	–	–	–	4.14	–	4.27	–	–	–	3.74
38	–	–	–	–	–	–	–	–	4.32	–
39	–	–	–	–	–	–	–	–	–	–
40	–	–	–	–	–	–	–	4.65	–	–

Table S13. Thurstone-scaled conversion table of YASR to YSR equivalents.

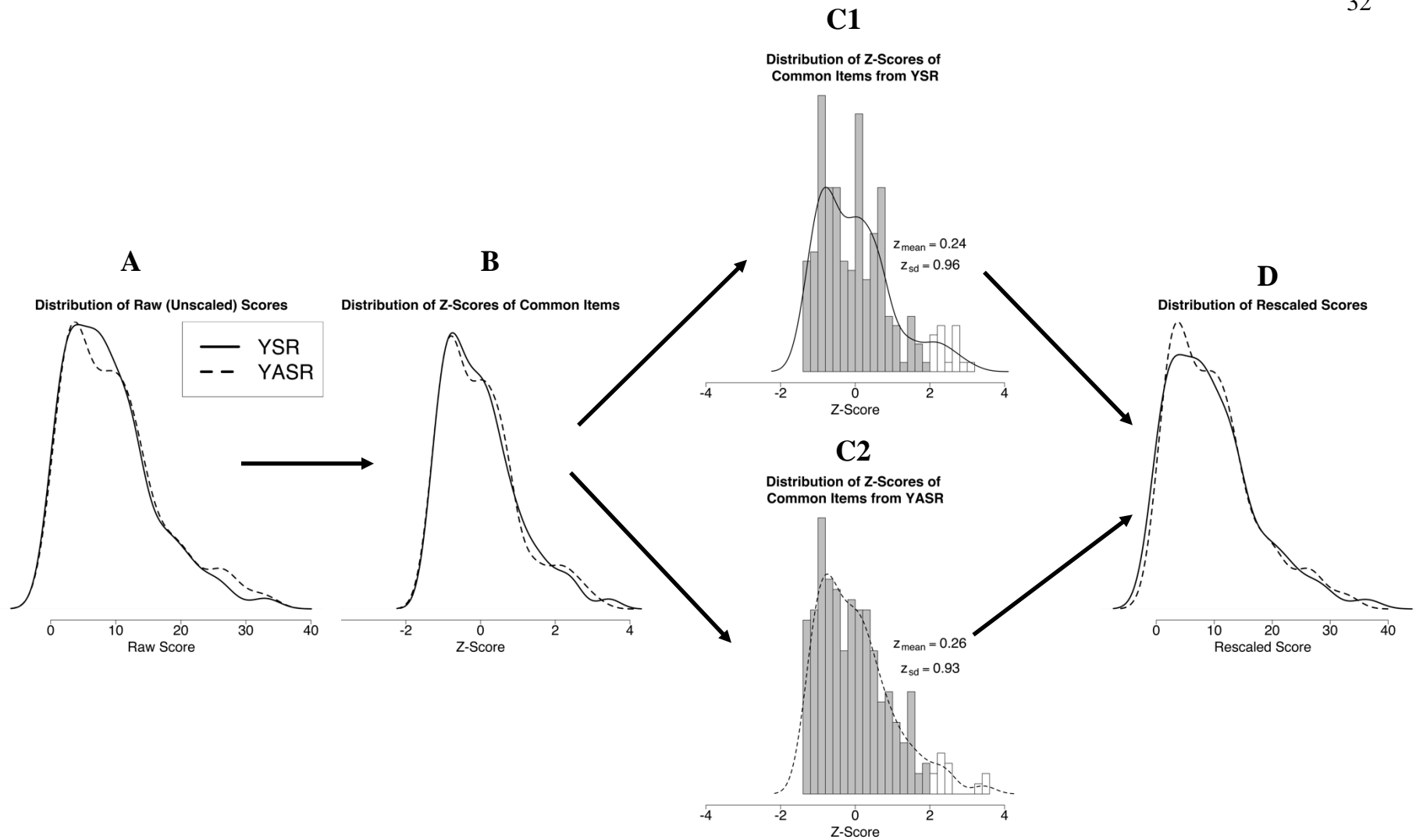
YASR score	Males					Females				
	Age (years)					Age (years)				
	20	21	22	23	24	20	21	22	23	24
0	-0.3	-0.4	-0.5	-0.4	-0.2	-0.6	-0.6	-0.5	-0.7	-0.6
1	0.8	0.6	0.5	0.6	0.8	0.5	0.5	0.6	0.4	0.5
2	1.8	1.6	1.5	1.6	1.8	1.6	1.6	1.7	1.5	1.6
3	2.8	2.6	2.6	2.6	2.8	2.7	2.7	2.8	2.6	2.7
4	3.8	3.6	3.6	3.7	3.8	3.8	3.8	3.9	3.7	3.8
5	4.8	4.7	4.6	4.7	4.8	4.9	4.9	5.0	4.8	4.9
6	5.8	5.7	5.6	5.7	5.8	6.1	6.0	6.1	5.9	6.0
7	6.8	6.7	6.7	6.7	6.8	7.2	7.1	7.2	7.0	7.1
8	7.8	7.7	7.7	7.7	7.8	8.3	8.2	8.3	8.1	8.2
9	8.8	8.8	8.7	8.8	8.8	9.4	9.3	9.4	9.2	9.3
10	9.9	9.8	9.7	9.8	9.8	10.5	10.4	10.5	10.3	10.5
11	10.9	10.8	10.8	10.8	10.8	11.7	11.5	11.6	11.4	11.6
12	11.9	11.8	11.8	11.8	11.8	12.8	12.6	12.7	12.5	12.7
13	12.9	12.9	12.8	12.8	12.8	13.9	13.6	13.8	13.6	13.8
14	13.9	13.9	13.8	13.9	13.8	15.0	14.7	14.9	14.7	14.9
15	14.9	14.9	14.9	14.9	14.8	16.1	15.8	16.0	15.8	16.0
16	15.9	15.9	15.9	15.9	15.8	17.2	16.9	17.2	16.9	17.1
17	16.9	16.9	16.9	16.9	16.8	18.4	18.0	18.3	18.0	18.2
18	17.9	18.0	17.9	17.9	17.8	19.5	19.1	19.4	19.1	19.3
19	19.0	19.0	19.0	18.9	18.8	20.6	20.2	20.5	20.2	20.4
20	20.0	20.0	20.0	20.0	19.8	21.7	21.3	21.6	21.3	21.5
21	21.0	21.0	21.0	21.0	–	22.8	22.4	22.7	22.4	22.6
22	22.0	22.1	–	22.0	21.8	24.0	23.5	23.8	23.5	23.7
23	23.0	23.1	–	23.0	22.8	25.1	24.6	24.9	24.6	24.8
24	24.0	24.1	24.1	24.0	23.8	26.2	–	26.0	25.7	25.9
25	25.0	25.1	25.1	25.1	24.8	27.3	–	27.1	26.8	27.0
26	26.0	26.2	26.1	–	25.8	28.4	27.9	28.2	27.9	28.1
27	27.0	–	–	–	26.8	29.5	29.0	–	29.0	29.2
28	–	28.2	28.2	28.1	27.8	30.7	–	30.4	–	30.3
29	–	–	–	–	28.8	–	31.2	31.5	–	–
30	–	–	–	–	29.8	–	32.3	–	–	32.5
31	–	–	–	–	–	–	–	–	33.3	–
32	–	–	–	–	–	35.1	–	–	–	34.7
33	–	–	–	–	–	36.2	35.5	–	–	–
34	–	–	–	–	33.8	37.4	–	37.1	36.6	36.9
35	35.1	–	–	–	–	–	37.7	–	–	38.0
36	–	–	–	–	–	–	38.8	–	–	–
37	37.2	–	–	–	36.8	–	–	40.4	–	–
38	–	–	–	38.3	–	–	–	–	–	–
39	–	–	–	–	–	–	–	–	–	–
40	–	–	40.5	–	–	–	–	–	–	–
<i>M</i>	7.77	7.68	7.79	7.78	8.50	9.50	9.73	10.26	10.12	10.53
<i>SD</i>	6.88	6.88	7.03	7.07	7.55	7.63	7.80	7.64	7.79	8.17

Note: values reflect the YASR scores on the scale of the YSR at age 19. Mean and standard deviation reflect the mean and standard deviation of the participants' re-scaled YASR scores on the YSR scale at age 19 (they do not reflect the mean and standard deviation of the values in the above column).

Table S14. Linear growth curve model of Thurstone-scaled internalizing problems.

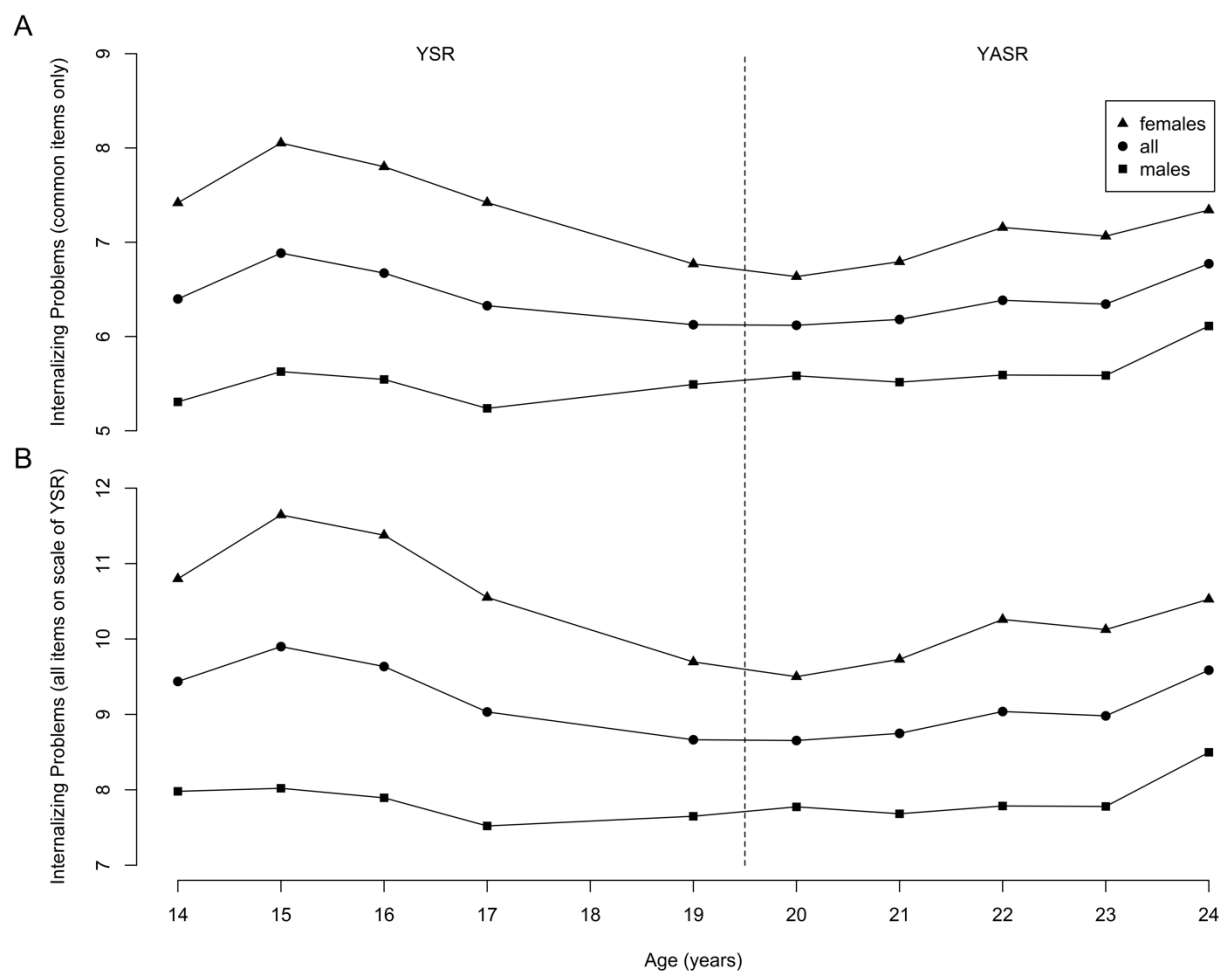
Variable	<i>B</i>	$\beta$	<i>SE</i>	<i>DF</i>	<i>p</i>
intercept	8.225	0.015	0.445	3977	< .001
time	-0.024	-0.033	0.057	3977	.674
Predictors of the intercepts					
female	3.267	0.181	0.602	539	< .001
African American	-1.331	-0.065	0.829	539	.109
Other Ethnicity	-2.610	-0.028	2.540	539	.305
Predictors of the slopes					
female	-0.102	-0.022	0.077	3977	.187
African American	-0.011	-0.002	0.110	3977	.920
Other Ethnicity	0.162	0.009	0.322	3977	.614
Variance components					
	<i>SD</i>				
intercept	6.19				
time	0.74				
residual	4.26				
Correlation between intercept and slope		<i>r</i> = .47			
Model Pseudo- <i>R</i> <sup>2</sup>		.747			



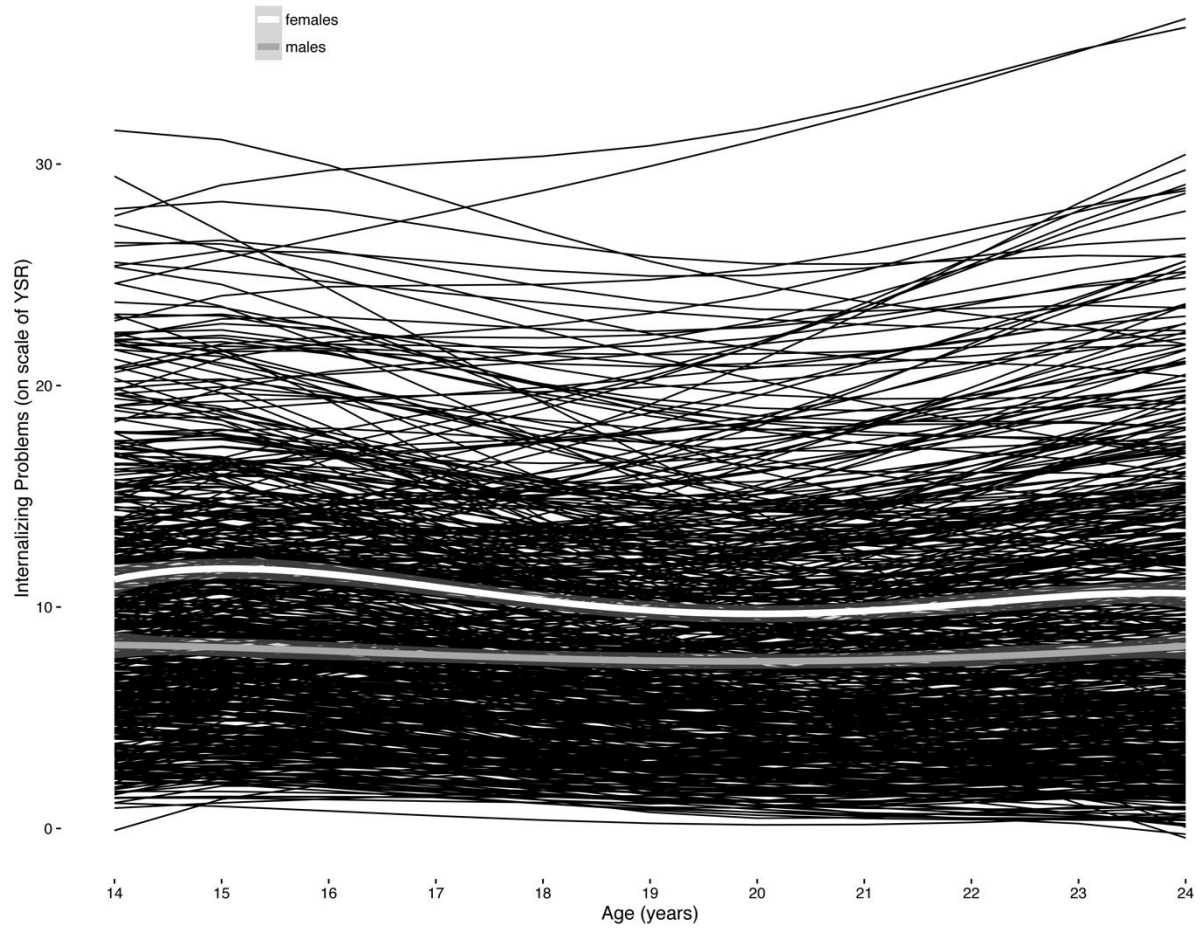


*Figure S1.* Depiction of steps in vertical scaling using Thurstone scaling with a common-item design. YSR = Youth Self-Report at age 19 (target scale). YASR = Young Adult Self-Report at age 20. Panel A depicts the raw score distributions of the two measures (distributions are depicted with kernel density estimation). Panel B depicts the distribution of  $z$ -scores of the items that are common to

both measures (i.e., the common items). Panel C depicts the distribution of  $z$ -scores of the common items for each measure (Panel C1 = YSR, Panel C2 = YASR), along with the calculations of the mean and standard deviation of  $z$ -scores within the target range of -2 to +2. Each histogram bar reflects the frequency of a given  $z$ -score (corresponding to a given raw score) on the measure. Gray histogram bars reflect  $z$ -scores within the target range of -2 to +2 that were used for calculating the mean and standard deviation. Note that the  $z$ -score for each unique raw score i.e., gray histogram bar, is used in the calculation (rather than all observed  $z$ -scores), so the mean and standard deviation do not necessarily equal 0 and 1, respectively. The measures are rescaled to be on the same scale by using the mean and the standard deviation of the  $z$ -scores of the common items to align their percentile scores. Panel D depicts the rescaled scores (i.e., scores from the YASR on the scale of the YSR). The mean and standard deviation of the rescaled scores were calculated using Equations S1 and S2, respectively. We calculated a conversion table by multiplying the  $z$ -scores of the total raw scores by the standard deviation of the scaled score and added the mean of the scaled score (see Table 2). The figure shows that, in comparison to the YSR, the unscaled YASR scores were over-represented at lower levels of the scale and under-represented at upper levels of the scale (presumably because of fewer items in the YASR; see Panel A). Rescaling the scores made the scales more comparable. Note that, by design, the distributions of rescaled scores for the two measures do not perfectly overlap. Vertical scaling does not create the same distribution (mean and standard deviation) for each measure because it retains differences in means and variances across the two measures (based on the means and variances of the common items). Nevertheless, the scores are on a more comparable scale. Although the common items are used to determine the general form of change on the same scale, all developmentally relevant, construct-valid items are used to estimate each person's trait level on this scale.



*Figure S2.* Panel A depicts participants' mean scores on the *common* items (i.e., the items that were common to the Internalizing scale of the Youth Self-Report, YSR, and Young Adult Self-Report, YASR). Panel B depicts participants' mean internalizing problem scores on *all* age-relevant items of the Internalizing scale, after rescaling the YASR scores to the metric of the YSR (based on the scale of the YSR at age 19) using Thurstone scaling. Internalizing problems to the left of the dashed line (i.e., ages 14–19) were rated on the YSR. Internalizing problems to the right of the dashed line (i.e., ages 20–24) were rated on the YASR. Internalizing problem reports were not collected at age 18.



*Figure S3.* Individuals' fitted quartic trajectories of Thurstone-scaled internalizing problems in black. Average quartic trajectory for females in white. Average quartic trajectory for males in gray.

## References

- Chalmers, R. P. (2012). mirt: a multidimensional item response theory package for the R environment. *Journal of Statistical Software*, *48*, 1-29.
- Chen, W.-H., & Thissen, D. (1997). Local dependence indexes for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, *22*, 265-289.  
doi:10.3102/10769986022003265
- Fennessy, L. M. (1995). *The impact of local dependencies on various IRT outcomes*. (Doctoral dissertation), Available from ProQuest Dissertations & Theses database. (UMI No. 9524701)
- Grimm, K. J., Ram, N., & Hamagami, F. (2011). Nonlinear growth curves in developmental research. *Child Development*, *82*, 1357-1371. doi:10.1111/j.1467-8624.2011.01630.x
- Hankin, B. L., Abramson, L. Y., Moffitt, T. E., Silva, P. A., McGee, R., & Angell, K. E. (1998). Development of depression from preadolescence to young adulthood: Emerging gender differences in a 10-year longitudinal study. *Journal of Abnormal Psychology*, *107*, 128-140. doi:10.1037/0021-843x.107.1.128
- Knight, G. P., & Zerr, A. A. (2010). Informed theory and measurement equivalence in child development research. *Child Development Perspectives*, *4*, 25-30. doi:10.1111/j.1750-8606.2009.00112.x
- Kolen, M. J., & Brennan, R. L. (2014). *Test equating, scaling, and linking: Methods and practices* (3rd ed.). New York, NY, US: Springer.
- Markon, K. E., Chmielewski, M., & Miller, C. J. (2011). The reliability and validity of discrete and continuous measures of psychopathology: A quantitative review. *Psychological Bulletin*, *137*, 856-879. doi:10.1037/a0023678

- Meade, A. W. (2010). A taxonomy of effect size measures for the differential functioning of items and scales. *Journal of Applied Psychology, 95*, 728-743. doi:10.1037/a0018966
- Morizot, J., Ainsworth, A. T., & Reise, S. P. (2007). Toward modern psychometrics: Application of item response theory models in personality research. In R. W. Robins, R. C. Fraley, & R. F. Krueger (Eds.), *Handbook of research methods in personality psychology* (pp. 407–421). New York, NY, US: Guilford Press.
- Raju, N. S. (1990). Determining the significance of estimated signed and unsigned areas between two item response functions. *Applied Psychological Measurement, 14*, 197-207. doi:10.1177/014662169001400208
- Tomarken, A. J., & Waller, N. G. (2003). Potential problems with 'well fitting' models. *Journal of Abnormal Psychology, 112*, 578-598. doi:10.1037/0021-843X.112.4.578
- van der Ark, L. A. (2007). Mokken scale analysis in R. *2007, 20*, 19. doi:10.18637/jss.v020.i11