

# BMJ Open Strengths and Difficulties Questionnaire: internal validity and reliability for New Zealand preschoolers

Paula Kersten,<sup>1</sup> Alain C Vandal,<sup>2,3</sup> Hinemoa Elder,<sup>4</sup> Kathryn M McPherson<sup>5,6</sup>

**To cite:** Kersten P, Vandal AC, Elder H, *et al.* Strengths and Difficulties Questionnaire: internal validity and reliability for New Zealand preschoolers. *BMJ Open* 2018;**8**:e021551. doi:10.1136/bmjopen-2018-021551

► Prepublication history for this paper is available online. To view these files, please visit the journal online (<http://dx.doi.org/10.1136/bmjopen-2018-021551>).

Received 8 January 2018  
Revised 7 March 2018  
Accepted 27 March 2018



<sup>1</sup>School of Health Sciences, University of Brighton, Brighton, UK

<sup>2</sup>Department of Biostatistics and Epidemiology, AUT University, Auckland, New Zealand

<sup>3</sup>Health Intelligence and Informatics, Ko Awatea, Counties Manukau District Health Board, Auckland, New Zealand

<sup>4</sup>School of Graduate Studies, Te Whare Wānanga o Awanuiārangī, Auckland, New Zealand

<sup>5</sup>Health Research Council of New Zealand, Auckland, New Zealand

<sup>6</sup>Centre for Person Centred Research, School of Clinical Sciences, AUT University, Auckland, New Zealand

## Correspondence to

Professor Paula Kersten;  
p.kersten@brighton.ac.uk

## ABSTRACT

**Objectives** This observational study examines the internal construct validity, internal consistency and cross-informant reliability of the Strengths and Difficulties Questionnaire (SDQ) in a New Zealand preschool population across four ethnicity strata (New Zealand European, Māori, Pasifika, Asian).

**Design** Rasch analysis was employed to examine internal validity on a subsample of 1000 children. Internal consistency (n=29 075) and cross-informant reliability (n=17 006) were examined using correlations, intraclass correlation coefficients and Cronbach's alpha on the sample available for such analyses.

**Setting and participants** Data were used from a national SDQ database provided by the funder, pertaining to New Zealand domiciled children aged 4 and 5 and scored by their parents and teachers.

**Results** The five subscales do not fit the Rasch model (as indicated by the overall fit statistics), contain items that are biased (differential item functioning (DIF)) by key variables, suffer from a floor and ceiling effect and have unacceptable internal consistency. After dealing with DIF, the Total Difficulty scale does fit the Rasch model and has good internal consistency. Parent/teacher inter-rater reliability was unacceptably low for all subscales.

**Conclusion** The five SDQ subscales are not valid and not suitable for use in their own right in New Zealand. We have provided a conversion table for the Total Difficulty scale, which takes account of bias by ethnic group. Clinicians should use this conversion table in order to reconcile DIF by culture in final scores. It is advisable to use both parents and teachers' feedback when considering children's needs for referral of further assessment. Future work should examine whether validity is impacted by different language versions used in the same country.

## INTRODUCTION

Educational achievement and problems in primary and secondary school aged children can arise as a result of behavioural and emotional problems when the child is of preschool age.<sup>1–5</sup> Consequently, screening to identify children with or at risk of behavioural problems at a preschool age is an increasingly used preventative strategy, aiming to enhance the success of support programmes and early intervention.<sup>6</sup> Such screening is best

## Strengths and limitations of this study

- A key strength of this study is the inclusion of all 4-year-old and 5-year-old children in New Zealand for whom a Strengths and Difficulties Questionnaire assessment was available in 2011, resulting in our ability to assess the validity of the tool at the population level and with sufficient power to make sound conclusions.
- A strength of the study included robust data quality checks and the exclusion of 39% of cases for which we had concerns about their quality (it being incomplete or containing multiple inconsistencies).
- A limitation was our inability to assess differential item functioning by other key variables that may affect validity, for example, first language or country of birth, as such data were not available.
- Future work should examine whether validity is impacted by different language versions used (in the same country).

performed using standardised methods, and for behavioural assessment, this means the use of a questionnaire-based measure. The Strengths and Difficulties Questionnaire for parents (SDQ-P) and for teachers (SDQ-T) is a tool used worldwide for this purpose to screen preschool children's psychosocial attributes (positive and negative behaviours).<sup>7–10</sup> It consists of 25 items, making up five subscales: Emotional Symptoms, Conduct Problems, Hyperactivity, Peer Problems and Prosocial Behaviour.<sup>7 8</sup>

Before using a measure such as the SDQ, establishing validity and reliability is key for optimum decision-making. At present, there are two dominant approaches to the development and testing of measures: Classical Test Theory (CTT) and Modern Test Theory (also known as item response theory).<sup>11</sup> In CTT, it is assumed that the observed scores on items are the sum of the true score (which we cannot directly measure) and measurement error. However, neither the true score nor the measurement error can be determined and the approach is therefore flawed.<sup>12</sup> In

addition, the best conclusion that can be made following satisfactory tests of validity and reliability using CTT is that an outcome measure is an ordinal scale. Yet, many statistical tests that examine the validity of scales assume that the data arising are of interval nature. Indeed, in the preschool population, the SDQ has only been tested using parametric, CTT approaches, as demonstrated in our recent systematic review<sup>13</sup> to which we return below. By contrast, Modern Test Theory approaches, such as Rasch analysis, are underpinned by mathematical models that specify the conditions under which equal interval measurements can be estimated from outcome measurement data.<sup>14–16</sup> These approaches are therefore more robust.

Evaluations of the structural validity of the SDQ drawing on CTT in preschoolers has been extensively researched using factor analysis (eg, by Klein *et al*, Tobia *et al* and Mieloo *et al*<sup>17–19</sup>), Cronbach's alphas ( $\alpha$ )<sup>13</sup> and correlation coefficients<sup>13 20</sup> and Weighted Least Squares in older children.<sup>21</sup> Our systematic review found acceptable to good evidence for the 5-factor SDQ structure in preschoolers, when confirmatory factor analysis (CFA) had been used.<sup>13</sup> A different approach to examining structural validity, using Modern Test Theory, can be achieved by examining whether each of the subscales are unidimensional and fit the Rasch model (ie, examining internal construct validity).<sup>15</sup> Like CFA, Rasch analysis is a confirmatory approach to examining whether items belong to the subscales under investigation. However, there are known limitations to using factor analysis on ordinal scales, including its parametric basis and the emergence of 'difficulty factors', which may spuriously indicate multidimensionality.<sup>22</sup> In addition, factor analysis does not allow detailed investigation of item function in regard to targeting, differential item functioning (DIF) and local dependency between items, whereas Rasch analysis includes such assessments.<sup>23</sup> We identified one study which had employed Rasch analysis on SDQ data that had been self-completed by 12–18-year olds in Sweden.<sup>24</sup> This study showed that none of the SDQ scales was psychometrically robust, with misfitting items in all five subscales and poor internal consistency. However, that study did not examine whether the scale was invariant across different subgroups.

Internal consistency of the SDQ-P subscales has been reported in many studies and synthesised in a systematic review.<sup>13</sup> The sample size-weighted average Cronbach's  $\alpha$  for the five subscales was below the threshold of 0.70 (implying inadequate internal consistency for shorter, established scales) and for the Difficulty scale  $\alpha$  was 0.79 (acceptable for group comparisons but not for individual use) (Streiner and Norman, p. 91).<sup>25</sup>

Inter-rater reliability of SDQ subscales between two parents and between two teachers has previously been found to be acceptable when correlation coefficients were used (between 0.42 and 0.64 for parents and between 0.59 and 0.81 for teachers).<sup>20</sup> Other studies have examined scores between different types of informants (eg,

parent and teacher). The systematic review showed that the sample size-weighted average correlation coefficients generated from these studies were weak to moderate (between 0.25 and 0.45).<sup>13</sup>

The validity and reliability of the SDQ have not previously been examined in New Zealand, a country with a sizeable indigenous population (Māori, 15.4%) and immigrant population (25.2% born overseas).<sup>26</sup> New Zealand is a multicultural society, impacting on values, ways of living and languages spoken. It cannot be assumed that measures capturing psychological constructs will have cultural equivalence.<sup>27 28</sup> Indeed, a New Zealand qualitative study has shown that parents from Māori, Pacific Island, Asian and new immigrant groups questioned the cultural validity of the SDQ.<sup>29</sup> Cultural equivalence therefore needs further investigation.

In summary, the use of CTT approaches to examine the validity of the SDQ are limited, evidence suggests cross-informant reliability is weak and there is no evidence for cultural equivalence for the New Zealand population. Therefore, we aimed to use Modern Test Theory, and specifically Rasch analysis, to examine the internal construct validity and cultural equivalence of the SDQ in a New Zealand preschool population across different ethnicity strata and to examine reliability between parents and teachers (cross-informant reliability). We hypothesised that the SDQ subscales and the Difficulty scale would (1) have cross-informant reliability (with consistency in scores by parents and teachers); (2) fit the Rasch model (demonstrating unidimensionality and internal construct validity) and (3) have cultural equivalence across ethnic strata (demonstrated by an absence of DIF).

## METHODS

### Study design and sample

This observational study used SDQ data gathered during the New Zealand Before School Check (B4SC), which takes place when the child is aged (4 or exceptionally aged 5).<sup>9</sup> The B4SC is carried out by registered nurses based in primary care and involves the assessment of the child's general health, hearing, oral health, vision, growth as well as developmental and behavioural problems. The latter is evaluated using the Australian SDQ version for 2–4-year olds, completed by the parent. If the child is in preschool, the nurse also requests their teacher to complete the SDQ for the child. Clear instructions for the administration of the SDQ are provided within the B4SC handbook. In New Zealand, there is no other SDQ data collection point during childhood.

Data sources/quality, missing data and bias: Permission to use the full, deidentified 2011 national B4SC SDQ dataset for preschoolers (n=51 251) from the New Zealand Ministry of Health was provided by the B4SC Governance Board. Data quality checks on SDQ data resulted in the deletion of 20 024 cases (out of n=51 251, 39%) for the following reasons:

1. Individual item data from the parent questionnaire were missing completely ( $n=19\,197$ ) or partially ( $n=1$ ) since (1) we would not have been able to carry out a quality check of the subscale scores and (2) we would not be able to use these data for the Rasch analysis); thus, 19 198 were removed from the analysis set.
2. District Health Boards (DHB) for which we had fewer than 15% of data on individual items, since the quality of their data is in doubt: although a total of 12 720 records came from these DHBs, this extra step only entailed the removal of a further 375 records from the analysis set after step 1.
3. Children's ages were recorded as younger than 4 or older than 5 when the SDQ was completed (we suspect some of these ages may have been entered incorrectly; however, this step only entailed the removal of a further 451 records from the analysis set after steps 1 and 2.
4. Cases with all zero scores: these were deemed potentially erroneous as the Prosocial subscale is scored in the opposite direction from the other subscales; although 1038 cases fitted this profile, none had complete parental item data and so no further record was removed on the basis of this criterion after steps 1, 2 and 3.

Study size: In total, 29 075 cases remained in the parents' dataset; 17 006 remained for the parent-teacher cross-informant reliability analysis. Rasch analysis uses fit statistics, but these are not suited to such large sample sizes. Fit to the Rasch model is considered acceptable when the observed data fit the predetermined Rasch model,<sup>15 30</sup> traditionally examined with fit statistics (eg, the item-trait interaction  $\chi^2$ ). A non-significant  $\chi^2$  indicates fit to the Rasch model. Power increases with large samples, which inflates the  $\chi^2$  and results in negligible small differences appearing as a statistically significant misfit between the data and the model.<sup>31 32</sup> Therefore, our Rasch analysis was carried out on a smaller sample ( $n=1000$ ), to allow examination of convergence to the Rasch model. The sample was created by randomly sampling equal numbers of cases from the total parent sample, for four main ethnic groups (250/ethnic group): New Zealand European (NZE), Māori, Asian and Pasifika. This is well above the recommended sample size for studies using Rasch analysis. For example, it has been suggested that to have 99% confidence that the estimated item difficulty is within  $\pm\frac{1}{2}$  logit of its stable value on the interval metric, the minimum sample size range is 108–243 (best to poor targeting).<sup>33 34</sup>

## Instruments

The SDQ consists of 25 items, each with three response options: not true, somewhat true and certainly true. The four SDQ subscales reflecting problematic behaviours or emotions (Emotional Symptoms, Conduct Problems, Hyperactivity, Peer Problems) contain 15 positively worded items and 5 negatively worded items.<sup>7 8</sup> Positively worded items are reverse scored (in New Zealand this is done on data entry); thus, higher subscale scores denote

## Box 1 Calculation of root mean square error of approximation (RMSEA)

In Rasch analysis, RMSEA is calculated as follows:

$$RMSEA = \sqrt{[(\chi^2/df) - 1]/(N - 1)}, 0$$

$$RMSEA = \sqrt{[(\chi^2/df) - 1]/(N - 1)}, 0^{32}$$

$\chi^2$  is the item-trait interaction chi-square (obtained from the analysis within the Rasch software),  $df$  is its degrees of freedom.

$N$  is the sample size.

Notice that the RMSEA has an expected value of 0 when the data fit the model. Overfit of the data to the model,  $\chi^2/df < 1$ , is ignored. For a given  $\chi^2$ , RMSEA decreases as sample size ( $N$ ) increases.

greater problems. Scores from these four subscales are also summed to give an overall Difficulty score ranging from 0 to 40. The five items making up the Prosocial Behaviour subscale are positively worded and higher scores denote better social behaviour.

## Data analysis

Cross-informant reliability (between parents and teachers) was assessed for those cases for which both parent and teacher SDQ data were available ( $n=17\,006$ ). The intraclass correlation coefficient (ICC) is the preferred statistical technique and was used.<sup>25 35</sup> However, as many studies of the SDQ have used correlations,<sup>36</sup> we will also present those.

Each SDQ subscale and the Difficulty scale were fitted to the Rasch model to examine fit, using RUMM2030 software.<sup>37</sup> Fit was considered acceptable if there was a non-substantial deviation of individual items and respondents from the Rasch model (individual item and person fit residuals should be within the range of  $\pm 2.5$ , the average fit residual statistics should be close to a mean of 0 and SD of 1, the item  $\chi^2$  should be non-significant). In addition, we used the root mean square error of approximation (RMSEA) to examine fit, with  $RMSEA < 0.02$  suggesting data fit the Rasch model (box 1).<sup>32</sup>

Log-transformed item scores generated from the response choices should reflect the increasing or decreasing latent trait to be measured (threshold ordering).<sup>30</sup> When a given level of problems is not confirmed by the expected response option to an item, disordered thresholds are observed. Disordering is only considered statistically significant if the 95% CI of the threshold locations do not overlap. When significant disordering is observed, response categories can be combined.

An assumption of the Rasch model is that the answers to one item should not be dependent on the responses to another item, conditional on the trait being measured. This local independence is examined by exploring the correlations between items' residuals, which should not be more than 0.20 above the average residual correlation.<sup>38</sup> If locally dependent items are observed, they can

be combined into a testlet, a bundle of items that share a common stimulus.<sup>39</sup>

The Rasch model expects that each item is invariant (unbiased) across key groups (eg, ethnicity or gender),<sup>40,41</sup> examined statistically with an analysis of variance and visually by examining the item characteristic curves. Variance (DIF) can be uniform; the bias is present consistently across the trait. For example, uniform DIF by ethnic group implies that item difficulty is different for individual ethnic groups across the trait even though their underlying level of problems is the same. DIF can also be non-uniform; the bias is not consistent across the trait. DIF analysis is affected by large sample sizes with non-significant DIF showing as significant; hence, inspection of item characteristic curves is also important. When uniform DIF is observed, two strategies can be employed. First, DIF items (if present in >1 item) can be combined into a testlet to examine if DIF is cancelled out at the test level; second, the item can be split by the variable for which DIF is observed. In our analysis, we considered the final solution to be the one with the best improvements in fit statistics.

Another key assumption of the Rasch model is that a scale must be unidimensional. This is examined by creating two subsets of items, identified by a principal component analysis of the item residuals, with those loading negatively forming one set and those positively loading the second set.<sup>42</sup> An independent t-test is used to compare estimates derived from the two subtests for each respondent. When fewer than 5% of the t-tests are significant (or the 95% CI of t-tests includes 5%), unidimensionality is supported.<sup>42,43</sup>

Targeting of the subscales to the population was examined with person-item-threshold maps.

Internal consistency was examined with Cronbach's  $\alpha$  and Person Separation Index (PSI) statistics. PSI is an indicator of the number of statistically different strata (groups) that the test can identify in the sample.<sup>44</sup> Interpretation of the PSI is similar to Cronbach's  $\alpha$  with values  $\geq 0.70$  suitable for group comparisons and  $\geq 0.85$  for individual clinical use. However, Cronbach's  $\alpha$  can only be calculated when there are no missing data and is not considered robust with skewed data.<sup>45</sup> Therefore, we present PSI and Cronbach's  $\alpha$  in summary tables as well as the number of groups between which the subscale is able to discriminate.<sup>46</sup>

Finally, for polytomous scales, two Rasch models can be used. The Rating Scale version assumes that the distance between thresholds is equal across items.<sup>14</sup> The Unrestricted (Partial Credit) model does not make this assumption.<sup>47</sup> A log-likelihood test examines whether results from these two models are significantly different and if this is so the Partial Credit model should be used. This test was significant ( $p < 0.001$ ) for all subscales and therefore the Partial Credit model was used.

### Patient and public involvement

End users of our research include families, preschool teachers, service providers and the Ministry of Health.

The research aims and questions were part of a tender prepared by the Ministry of Health, to which we responded. Thus, we did not have the ability to include end users in the development of study questions. The analysis presented here did not require participant recruitment or data collection and end users were therefore not consulted about the study design. Researchers in New Zealand have a responsibility to ensure their research is of value and culturally responsive to Māori. Therefore, guidance for the study was sought from the University's Mātauranga Māori committee, which members are drawn from a wide range of Māori communities. The findings from the part of the study reported here were presented to the Ministry of Health.

### RESULTS

The child gender split was balanced with 49% female and 51% male in the full parent sample as well as the cross-comparison sample; 99.6% were aged 4 at the time of the B4SC (0.4% of children had recently turned 5). Child ethnicity in the parent sample was 57% NZE, 23% Māori, 12% Pasifika and 8% Asian; this distribution was similar in the cross-comparison sample 63% NZE, 16% Māori, 7% Pasifika and 7% Asian. As noted above, there were no missing data in the selected samples.

#### Cross-informant reliability (n=17 006)

Cross-informant reliability between parent and teachers as measured by correlations was generally poor (all  $< 0.5$ , mean 0.28) and ICCs (all  $< 0.6$ , mean 0.13). Cross-informant reliability was better in the Hyperactivity subscale and worst in the Prosocial subscale, better for NZE and worst for Pasifika children (table 1).

#### Internal validity and cross-cultural equivalence

Table 2 displays results from the Rasch analysis.

#### Emotional Symptoms subscale

All items in this subscale had ordered thresholds, items were locally independent and the subscale was unidimensional. Person fit was adequate with a mean person fit residual reasonably close to 0 and the SD below 1.4 (table 2: analysis 1). However, overall fit to the Rasch model was unsatisfactory (RMSEA  $> 0.02$ ). PSI was below 0 and Cronbach's  $\alpha$  0.15. All item fit residuals were within the acceptable range of  $-2.5$  to  $2.5$ ; however, four out of five item  $\chi^2$  values were statistically significant, indicating misfit.

There was statistically significant uniform DIF by ethnicity in items 16 and 24, which was confirmed by visual inspection of the item characteristic curves (figure 1). Items 16 and 24 were combined into a testlet. This resulted in poorer person fit and similar RMSEA values (0.072). We therefore split these items by ethnic groups instead, creating unique items for NZE, Māori, Asian and Pasifika peoples, resulting in 11 items for the subscale. This step improved overall fit to the Rasch

**Table 1** Intraclass correlation coefficients SDQ subscales, overall and by ethnicity (n=17 006)

Variable	Ethnicity				
	Overall*	Māori	NZ European	Pasifika	Asian
	r	r	r	r	r
Valid N	17 056	2677	10 735	1144	1169
Mean item correlations	0.282	0.237	0.315	0.130	0.210
Minimum item correlations	0.199	0.151	0.220	-0.009	0.055
Maximum item correlations	0.418	0.358	0.447	0.275	0.377
	ICC	ICC	ICC	ICC	ICC
Emotional Symptoms	0.126	0.067	0.186	0.017	0.098
Conduct Problems	0.137	0.112	0.179	0.038	0.079
Hyperactivity	0.174	0.136	0.245	0.050	0.122
Peer Problems	0.139	0.100	0.202	0.004	0.162
Prosocial	0.055	0.048	0.066	0.040	0.035
Mean ICC	0.126	0.093	0.175	0.030	0.099
Minimum ICC	0.055	0.048	0.066	0.004	0.035
Maximum ICC	0.174	0.136	0.245	0.050	0.162

SDQ, Strengths and Difficulties Questionnaire.

\* Overall: all four ethnic groups combined

ICC: intraclass correlation coefficient

model; however, the RMSEA was still greater than the acceptable value of 0.02 and internal consistency unacceptably low (table 2: analysis 2).

After items were split, all item fit residuals were within range, although two still had statistically significant  $\chi^2$  values (items 24NZE and item 8). Table 3 shows that the easiest item to endorse is item 16 and the hardest to endorse is item 13. The split item locations show that for children with the same level of Emotional Problems, item 16 is more readily endorsed when they are Māori and less readily endorsed when they are Pasifika (difference of 0.42 logits). Item 24 is endorsed more readily by parents of Asian than NZE children (difference of 0.49 logits). Figure 2 displays the targeting of the subscale to the population, clearly demonstrating the large number of extreme cases.

### Conduct Problems subscale

Conduct Problems item thresholds were ordered, items were locally independent and person fit and unidimensionality were acceptable. However, overall fit to the model was unsatisfactory (RMSEA>0.02, table 2: analysis 3). Internal consistency was poor (PSI 0.10,  $\alpha$  0.65) with the subscale being able to discriminate between three strata.

Item fit residuals were within acceptable range though two had significant  $\chi^2$  (items 5 and 18).

Statistically significant DIF by ethnicity was present for item 12 and by gender for item 7. These two items were split by ethnicity and gender, respectively (table 2: analysis 4), resulting in satisfactory fit residuals, one item with a significant  $\chi^2$ , significant improvement in RMSEA

(0.03) but poor internal consistency (PSI=0.11, splitting items leads to missing data and  $\alpha$  cannot be calculated).

The easiest item to endorse was item 5 and the hardest item 12 (table 3). The split item locations show that for children with the same level of Conduct Problems, item 12 is more readily endorsed when they are Pasifika and less readily endorsed when they are NZE (difference of 1.22 logits). Item 7 is endorsed more readily by parents of boys than girls (difference of 0.32 logits). Targeting showed a floor effect (figure 2).

### Hyperactivity subscale

Ordered thresholds, local independence, person fit and unidimensionality were observed for the Hyperactivity subscale; however, overall fit to the model and internal consistency was unsatisfactory (RMSEA>0.02; PSI 0.30,  $\alpha$  0.48; subscale discriminates between three strata, table 2: analysis 5). Item fit residuals were out of range for item 21 and item 25 had a significant  $\chi^2$ . Uniform DIF was statistically significant by ethnicity in two items (15 and 21). These items were therefore split by ethnicity. This improved fit to the Rasch model (table 2: analysis 6) and displayed better fit than when these two items were combined into a testlet. Item fit residuals were within acceptable range of -2.5/+2.5; only one item had a significant item  $\chi^2$  statistic (table 3), and RMSEA was close to 0.02. However, internal consistency remained poor (PSI=0.31). The easiest item to endorse was item 15 (for Asian children) and the hardest item 10. The split item locations show that, for children with the same level of hyperactivity problems, item 15 is more readily endorsed when they are Asian and less readily endorsed when they

**Table 2** Fit to the Rasch model—SDQ-P (n=1000)

Subscales	Item fit residual				Person fit residual				$\chi^2$ interaction			Internal consistency†		Unidimensionality t-tests (CI)‡
	N	Mean§	SD	Mean	SD	Value	df	P value	RMSEA*	PSI without extremes	$\alpha$ without extremes	% (95% CI)		
<b>Emotional Symptoms</b>														
1 Initial	1000	-0.791	0.894	-0.327	0.783	83.6	20	<0.0001	0.068	-0.40	0.15	0		
2 Split items 16 and 24	1000	-0.545	0.841	-0.343	0.735	99.1	41	<0.0001	0.045	-0.41	-	0		
<b>Conduct Problems</b>														
3 Initial	1000	0.266	1.273	-0.253	0.876	71.6	20	<0.0001	0.060	0.10	0.65	0		
4 Split items 7 and 12	1000	0.134	0.902	-0.254	0.882	75.3	45	0.003	0.031	0.11	-	0		
<b>Hyperactivity</b>														
5 Initial	1000	0.260	2.348	-0.359	1.147	97.3	25	<0.0001	0.06	0.30	0.48	0.5 (-1.0 to 2.0)		
6 Split items 15 and 21	1000	0.323	1.480	-0.365	1.134	125.6	69	<0.0001	0.03	0.31	-	0.5 (-1.0 to 2.0)		
<b>Peer Problems</b>														
7 Initial	1000	-0.339	0.868	-0.207	0.719	69.0	20	<0.0001	0.06	-0.49	0.51	0		
8 Split item 23	1000	-0.207	0.652	-0.213	0.733	79.5	52	0.008	0.03	-0.43	-	0		
<b>Prosocial</b>														
9 Initial	1000	-0.075	1.592	-0.319	1.079	66.6	20	<0.0001	0.06	-0.03	0.29	0.1 (-1.5 to 1.8)		
<b>Difficulty</b>														
10 Initial	1000	-0.448	1.848	-0.248	1.004	296.3	180	0.0001	0.03	0.71	0.79	5.9 (4.6 to 7.3)		
11 Testlets DIF items and LD items	1000	-0.615	1.321	-0.294	0.985	200.4	144	0.001	0.02	0.71	0.77	3.0 (1.6 to 4.4)		

Indices indicative of fit:

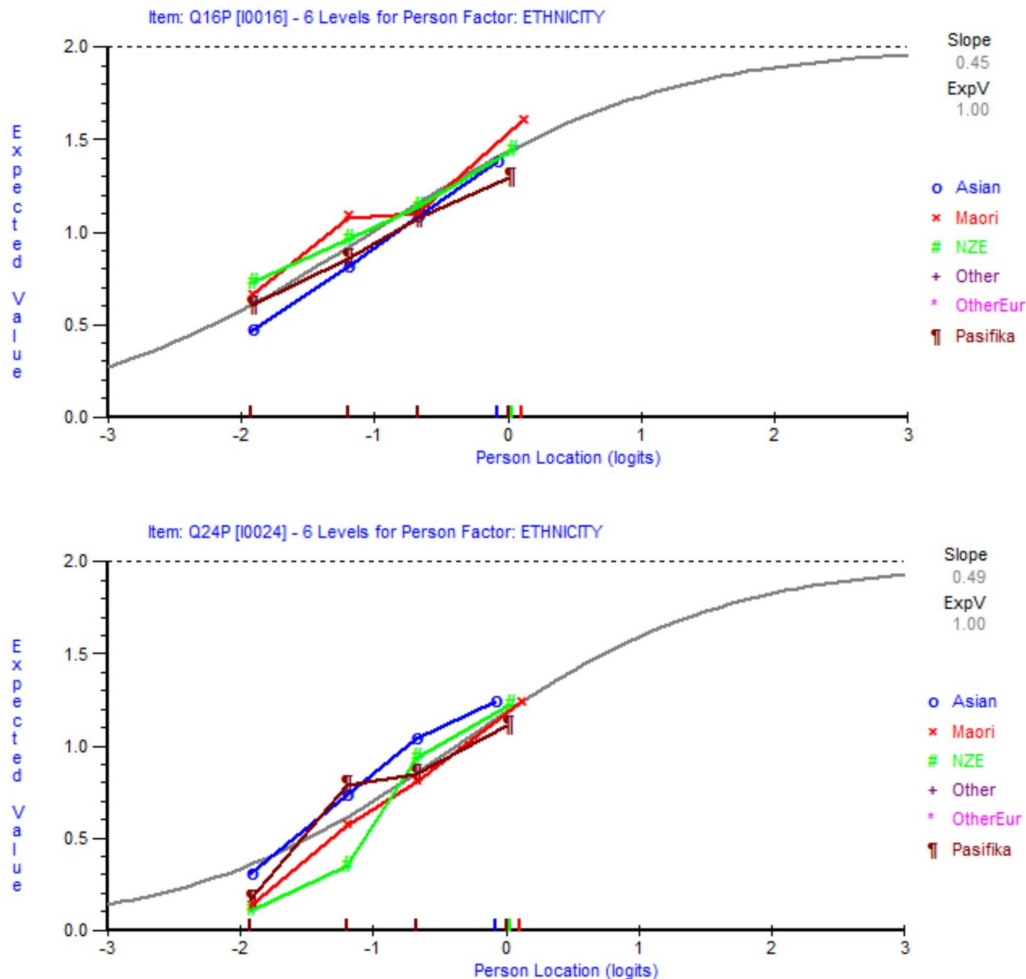
\*RMSEA<0.02.

†Internal consistency PSI and  $\alpha$ >0.70 (allows for group comparisons) and  $\geq$ 0.85 (allows for individual clinical use).

‡Unidimensionality indicated if fewer than 5% of t-tests are significant (ie, the 95% CI should include 5%).

§Mean item and person fit residuals: should be close to 0 (and<0.4); SD close to 1 (and <1.4).

DIF, differential item functioning; LD, local dependency; PSI, Person Separation Index; RMSEA, root mean square error of approximation; SDQ-P, Strengths and Difficulties Questionnaire parents.



**Figure 1** Item characteristics curves for items from the Strengths and Difficulties Questionnaire (parents, n=1000). NZE, New Zealand European.

are NZE (difference of 0.52 logits). Item 21 is endorsed more readily by parents of NZE children than Pasifika children (difference of 0.47 logits, [table 3](#)). The targeting map showed a floor effect ([figure 2](#)).

#### Peer Problems subscale

Ordered thresholds, local independence, person fit and unidimensionality were observed. However, overall fit to the Rasch model and internal consistency were unsatisfactory (RMSEA>0.02; PSI negative value,  $\alpha$  0.51, the subscale is able to discriminate between two strata, [table 2](#): analysis 7). Item fit residuals were acceptable, although two items had significant  $\chi^2$ . One item (23) displayed uniform DIF by ethnicity. After splitting this item by ethnicity, fit improved; all item fit residuals were within range (item 14  $\chi^2$  was borderline statistically significant), RMSEA was close to 0.02. PSI values remained negative however ([table 2](#): analysis 8). The easiest item was item 23 (for Asian children) and the hardest item 14. Item 23 was easier for Asian children and hardest for NZE children (difference of 1.10 logits, [table 3](#)). Targeting showed a significant floor effect ([figure 2](#)).

#### Prosocial subscale

The subscale met the requirements for threshold ordering, local independence, person fit and unidimensionality. Overall fit to the Rasch model and internal consistency were unsatisfactory (RMSEA>0.02; PSI negative values,  $\alpha$  0.29, subscale able to discriminate between two strata, [table 2](#): analysis 9). Item fit residuals were within the  $-2.5/+2.5$  range, though two had significant item  $\chi^2$  statistics. There was no DIF. Item 17 was the easiest to endorse; item 4 was the hardest to endorse. A ceiling effect was observed in the person-item-threshold map ([figure 2](#)).

#### Difficulty scale

Two items had disordered thresholds; however, this was not statistically significant and item response categories did not need to be combined. Some local dependency was present in two item pairs. Unidimensionality was observed ([table 2](#): analysis 10). Five item fit residuals were out of the acceptable range of  $-2.5/+2.5$  and four items showed uniform DIF by ethnicity (items 12, 16, 21 and 23). To examine whether DIF was present at the test level, these items were combined into a testlet. This

**Table 3** Item locations (in location order) and fit statistics SDQ-P subscales (n=1000)

Subscale and items	Location	SE	Fit residual	$\chi^2$ value	df	P value
<b>Emotional Problems*</b>						
16 Māori	-0.871	0.113	-0.226	2.968	4	0.5631
16 NZE	-0.692	0.124	-0.036	0.60	3	0.8960
16 Asian	-0.538	0.118	-0.101	0.77	3	0.8569
16 Pasifika	-0.450	0.120	0.911	0.61	3	0.8936
24 Asian	-0.250	0.124	-0.185	5.13	3	0.1629
24 Māori	0.010	0.117	-0.737	9.69	4	0.0461
24 Pasifika	0.024	0.124	-0.002	11.857	3	0.0079
24 NZE	0.243	0.127	-1.610	14.095	3	0.0028
3	0.653	0.070	-0.615	15.156	5	0.0097
8	0.908	0.075	-1.970	21.479	5	0.0007
13	0.965	0.080	-1.423	16.749	5	0.0050
<b>Conduct Problems†</b>						
5	-0.985	0.063	0.011	15.38	5	0.0089
18	-0.707	0.066	-0.352	22.19	5	0.0005
7 Male	-0.594	0.096	1.209	7.71	5	0.1732
7 Female	-0.271	0.100	1.917	6.09	5	0.2975
22	-0.012	0.072	0.156	8.49	5	0.1312
12 Pasifika	0.089	0.143	-0.148	3.527	5	0.6193
12 Māori	0.339	0.145	-0.512	5.862	5	0.3199
12 Asian	0.838	0.202	-0.030	2.344	5	0.7998
12 NZE	1.304	0.211	-1.049	3.733	5	0.5884
<b>Hyperactivity*</b>						
15 Asian	-0.491	0.109	-0.395	8.25	5	0.1432
15 Māori	-0.315	0.117	0.433	1.78	6	0.9388
21 NZE	-0.234	0.142	2.204	17.50	5	0.0037
2	-0.206	0.056	-1.327	23.29	9	0.0056
21 Asian	-0.186	0.124	1.414	8.216	5	0.1447
15 Pasifika	-0.019	0.121	0.388	8.775	5	0.1184
15 NZE	0.032	0.126	-1.737	12.772	5	0.0256
21 Māori	0.114	0.129	1.743	7.403	6	0.2852
21 Pasifika	0.234	0.122	1.393	5.986	5	0.3076
25	0.360	0.066	1.421	9.335	9	0.4070
10	0.712	0.065	-1.984	22.26	9	0.0081
<b>Peer Problems†</b>						
23 A	-0.968	0.109	-0.571	1.959	4	0.7432
23 P	-0.870	0.107	0.307	4.311	5	0.5056
23 M	-0.217	0.119	0.038	5.529	4	0.2372
6	-0.026	0.065	0.526	10.572	9	0.3062
23N	0.130	0.154	0.093	3.548	3	0.3147
11	0.233	0.066	-1.419	17.787	9	0.0377
19	0.491	0.071	0.131	12.305	9	0.1967
14	1.227	0.084	-0.763	23.501	9	0.0052
<b>Prosocial‡</b>						

Continued



Table 3 Continued

Subscale and items	Location	SE	Fit residual	$\chi^2$ value	df	P value
1	-0.487	0.079	-1.530	18.205	4	0.0011
4	-0.036	0.073	-0.273	12.624	4	0.0133
9	0.000	0.072	1.092	6.74	4	0.1502
17	0.008	0.071	-1.633	21.52	4	0.0003
20	0.515	0.073	1.972	7.52	4	0.1109
Difficulty§						
15	-0.835	0.054	-2.777	27.39	9	0.0012
LD items¶	-0.606	0.037	-1.744	14.01	9	0.1221
5	-0.583	0.056	-0.595	8.71	9	0.4645
DIF items**	-0.375	0.031	-2.500	21.03	9	0.0125
25	-0.331	0.061	0.036	14.05	9	0.1207
24	-0.314	0.058	0.839	7.44	9	0.5911
18	-0.313	0.059	-0.742	6.83	9	0.6553
6	-0.137	0.061	1.137	4.47	9	0.8777
7	-0.026	0.063	-1.305	23.26	9	0.0057
11	0.117	0.067	0.862	9.76	9	0.3702
22	0.308	0.068	-1.218	14.07	9	0.1199
3	0.311	0.071	1.017	11.50	9	0.2433
19	0.413	0.072	-1.247	10.59	9	0.3048
8	0.561	0.077	0.105	4.79	9	0.8525
13	0.646	0.087	0.621	9.37	9	0.4035
14	1.164	0.084	-2.326	13.15	9	0.1560

\*Bonferroni corrections applied p value is statistically significant if <0.005.

†Bonferroni corrections applied p value is statistically significant if <0.006.

‡Bonferroni corrections applied p value is statistically significant if <0.01.

§Bonferroni corrections applied p value is statistically significant if <0.003.

¶LD items; combined into a testlet (items 2 and 10).

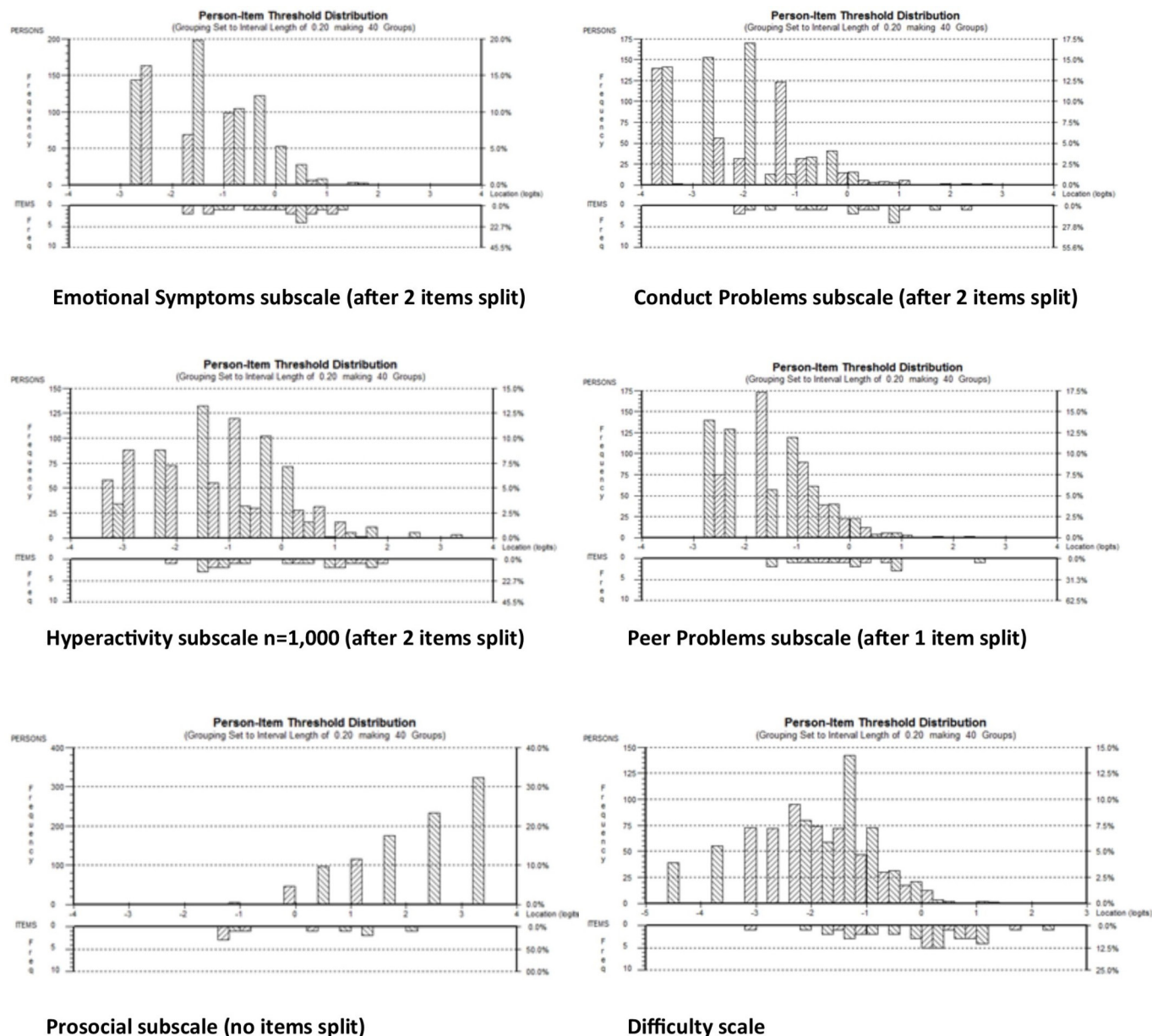
\*\*DIF items combined into a testlet (items 12, 16, 21, 23).

DIF, differential item functioning; LD, local dependency; NZE, New Zealand European; SDQ-P, Strengths and Difficulties Questionnaire parents.

resulted in an absence of DIF; however, one item pair remained locally dependent (items 2 and 10). A second testlet was created to deal with this local dependency. The resulting scale was unidimensional, with locally independent items (table 2: analysis 11). The RMSEA was within range suggesting overall fit to the Rasch model. Internal consistency was good (PSI 0.71,  $\alpha$  0.77, the scale was able to discriminate between six distinct strata). The fit residual for one item was slightly out of range (item 15, -2.777); however, given the negative value of this residual, this indicates redundancy rather than misfit and the item was therefore retained. The easiest item to endorse was item 15, the hardest item 14. The person-item threshold map showed a normal distribution, although located to the left of the item locations on the latent trait. A conversion table was produced, which can be used to convert the raw ordinal score to an interval scale (table 4).

## DISCUSSION

This study has shown that the SDQ items response categories work well; however, the five subscales diverge significantly from the Rasch model and four SDQ subscales include items that are biased by key variables with ethnicity having the greatest contribution. This raises critical questions about cultural equivalence. The five subscales suffer from a floor and ceiling effect and their internal consistency statistics are well below the acceptable range. By contrast, the Total Difficulty scale, which combines the four subscales capturing children's problems, is unidimensional, fits the Rasch model (after dealing with DIF and local dependency) and has internal consistency sufficient to distinguish between six groups of children. The study has also shown that parents and teachers score children in their care differently. Thus, all three study hypotheses are rejected. This section will discuss our findings in terms of fit to the Rasch model, internal consistency, cultural equivalence and cross-informant reliability.



**Figure 2** Person-item-threshold maps Strengths and Difficulties Questionnaire (parents, n=1000).

### Fit to the Rasch model

The Total Difficulty scale did fit the Rasch model, after dealing with four DIF items and two locally dependent items. This scale has good internal consistency and is able to discriminate between six groups of children on the latent trait. We observed the population distribution, while following a normal pattern, was to the left of the item locations on the latent trait. Thus, the precision of person estimates at the lower of the scale will not be as good as for those at the higher end of the scale. However, the SDQ is used for screening and arguably precise measurement at the lower end is not needed, since all one needs to establish is that the child does not need to be referred for further assessment or intervention. As we achieved fit to the Rasch model, we were able to provide a conversion table which can be used by clinicians to

convert the raw ordinal score to more accurate interval level and which takes account of DIF.

### Internal consistency

The five subscales are relatively short, which affects internal consistency and the subscales' ability to make fine distinctions between groups of people on the underlying trait.<sup>25</sup> In addition, there was significant divergence between the PSI and Cronbach's  $\alpha$  statistics, with PSI being much smaller than alpha. This divergence can be explained by the way these statistics are calculated. The calculation of Cronbach's  $\alpha$  assumes all SEs for individuals are the same, making it not a very robust statistics for skewed data.<sup>45</sup> This assumption results in relatively high values even in the presence of extreme scores and the Cronbach's  $\alpha$  values are therefore meaningless for SDQ

**Table 4** Conversion table for the Difficulty scale of the SDQ-P

Original Total Difficulty score (ordinal data)	Logit scores (interval level data)	Converted logit scores to 0–40 scale (interval level data)
0	-4.483	0
1	-3.655	4
2	-3.082	7
3	-2.685	8
4	-2.375	10
5	-2.117	11
6	-1.895	12
7	-1.699	13
8	-1.522	14
9	-1.36	15
10	-1.209	15
11	-1.068	16
12	-0.935	16
13	-0.809	17
14	-0.687	18
15	-0.571	18
16	-0.457	19
17	-0.347	19
18	-0.24	20
19	-0.134	20
20	-0.029	21
21	0.075	21
22	0.178	22
23	0.282	22
24	0.386	23
25	0.492	23
26	0.599	24
27	0.709	24
28	0.822	25
29	0.94	25
30	1.064	26
31	1.196	26
32	1.337	27
33	1.491	28
34	1.663	29
35	1.859	29
36	2.09	31
37	2.373	32
38	2.746	34
39	3.301	36
40	4.125	40

SDQ-P, Strengths and Difficulties Questionnaire parents.

data. This issue has not been raised in the SDQ literature; indeed, Cronbach's  $\alpha$  values are widely reported as satisfactory.<sup>48</sup> In Rasch analysis, the SE for every individual is estimated and the calculation of the PSI statistic takes these into account. Since SEs are largest for people with extreme scores, PSI will be smaller than Cronbach's  $\alpha$  as observed in our skewed data. However, the purpose of the SDQ is to identify those children who would benefit from further assessment or intervention. Thus, the fact that we observed a floor and ceiling effect is not necessarily problematic.

### Cultural equivalence

This study examined invariance by ethnicity at the item level and found lack of cultural equivalence. DIF (especially by ethnicity) was found for all the four subscales measuring problems, suggesting there are a number of questions to which parents respond differently despite overall scoring the same amount of problems on the trait being measured. The only other Rasch analysis study we were able to locate (conducted on data from children aged 12 to 18) did not include a DIF analysis and thus we cannot compare our findings against theirs.<sup>24</sup> Lack of measurement invariance of the subscales has also been shown by others (although on older children than in our sample) when using a CFA approach.<sup>50,51</sup> Richter *et al* found varying factor loadings and thresholds between different ethnic Norwegians and minority ethnic groups of adolescents and concluded that the total difficulty score is preferable.<sup>49</sup> Similarly, Ortuño-Sierra *et al* demonstrated that measurement variance was only partial, with 11 of the 25 items not being variant across different European samples.<sup>50</sup> By contrast, others have shown measurement invariance between British Indian and British white children using multigroup confirmatory factor analyses and demonstrated evidence of acceptable fit across ethnicity, although again their population was older (5–16 years) than the sample considered here.<sup>51</sup>

If measurement variance (DIF) is ignored, the child's difficulties can be overestimated or underestimated since the difficulty of the item varies by ethnic group, potentially leading to inaccurate identification of cases. This is important, given caseness has been shown to vary for different ethnic groups within the same country and between countries.<sup>52–54</sup> Our study is unable to assess why such DIF occurs, since the study drew on secondary data. However, we can pose some possible factors that may have affected measurement variance, as discussed below.

Our recent qualitative study suggests there is variation in the way the SDQ is administered—some parents complete the tool by themselves and others receive support from nurses, possibly impacting on the way questions are interpreted.<sup>29</sup> In addition, New Zealand preschool parents from Māori, Pacific Island, Asian and new immigrant groups questioned the cultural validity of the SDQ.<sup>29</sup> Respondents in an Australian qualitative study exploring the SDQ in Aboriginal community-controlled health services reported that the use of a questionnaire

as opposed to a general conversation or interview was deemed culturally inappropriate and that inter-relationships with peers were considered of less importance than relationships with family and participants.<sup>55</sup>

There are 85 different language versions available from the Youth in Mind website, though not one in Te Reo Māori (<http://www.sdqinfo.org/>). Translations and adaptations are not permitted without the involvement of that study team, which provides confidence in the robustness of translations. However, for our study, we do not know whether respondents were offered the SDQ in the language of their choice, as such data are not collected as part of the B4SC. The literature includes six studies that examined and demonstrated some issues with SDQ translations.<sup>13</sup> Using a language version that is not understood by respondents will affect validity,<sup>56</sup> which may have occurred here.

It is possible that poor literacy impacts on answering the SDQ, as found by others.<sup>57,58</sup> In New Zealand, there are many people (in proportion) with poorer than average literacy skills.<sup>59</sup> In addition, 18.6% of the New Zealand population report speaking two or more languages, the majority being born overseas (60.4%); many among these will have English as a second language.<sup>60</sup>

These aspects have particular relevance for Māori whānau (extended families) in New Zealand where it is estimated that 20% of Māori children and youth have Conduct Problems.<sup>61</sup> Therefore, it is important that screening of Māori children during the preschool years is accurate in ensuring that Māori whānau both receive the support they need and at the same time are not pathologised by false positive findings. The 2013 New Zealand Census found that 21% of the almost 700 000 Māori population could hold conversation about everyday things in Te Reo Māori, which has been a national official language since 1987.<sup>62</sup> Yet, there is not Māori version of the SDQ, or a New Zealand version incorporating commonly used Māori words.

### Cross-informant reliability

Cross-informant reliability was examined with ICCs which were well below the acceptable cut-off value of 0.6 (the mean in our study was 0.126). However, some argue that correlation coefficients can be used in the assessment of cross-informant reliability of the SDQ since parents and teachers make SDQ ratings based on different sources of information.<sup>7,48</sup> Our systematic literature review found weighted averages of coefficients between different informants ranged from 0.24 to 0.45,<sup>13</sup> similar to findings by others (range 0.26–0.47).<sup>48</sup> In our study, the mean correlation coefficient was 0.28, meaning only 8% of the variance can be explained by scores from different informants. This implies the importance of taking into account the views of both parents and teachers when making a decision for onward referral, a practice that is not commonplace in New Zealand.<sup>63</sup>

A key strength of this study is the inclusion of all preschool children in New Zealand for whom an SDQ

assessment was available in 2011, resulting in our ability to assess the validity of the tool at the population level, with sufficient power to make sound conclusions and ability to generalise to the wider New Zealand preschool population. Another strength was robust data quality checks and the exclusion of 39% of cases for which we had some concerns about quality (it being incomplete or containing multiple inconsistencies). From our steering group meetings, we gathered that there were a few reasons underlying these quality issues. In some DHBs, staff enter only the total scores, as opposed to item-level data. This practice leads to potential summing errors of total scores and these could not be checked or indeed analysed (hence we excluded these cases). Second, some DHBs told us they set the default values of answers as zero rather than blank. Consequently, when there were missing data (eg, if a teacher-completed SDQ was not available), the software would have summed these and arrived at total scores of 0. Given that the Prosocial scale is scored in the opposite direction of the others, zero scores on all subscales would be highly inconsistent and therefore shed doubt on data quality (and hence these were also excluded). An additional limitation was our inability to assess DIF by other key variables that may affect validity, for example, first language or country of birth, as such data were not available.

In conclusion, the Total Difficulty scale is internally valid and has acceptable internal consistency. Clinicians should use the conversion table as it accounts for bias by ethnic group. The five subscales are not valid and not suitable for use in their own right in New Zealand. Since consistency of scores between parents and teachers was poor, it is advisable to use both parents and teachers' feedback when considering children's needs for referral to further assessment. Future work should examine whether validity is affected by different language versions used (in the same country).

**Acknowledgements** We thank the funder for supporting the study.

**Contributors** PK conceived of the study, led on study design, project management, data analysis and dissemination. ACV, HE, KMMcP contributed to study design. ACV contributed to the data analysis. PK drafted the manuscript and is the guarantor. All authors revised it critically for important intellectual content and approved the final version for publication. All authors agree to be accountable for all aspects of the work.

**Funding** This work was supported by the Ministry of Health of New Zealand (grant number 341088).

**Disclaimer** All other authors declare no financial relationships with any organisations that might have an interest in the submitted work in the previous three years; no other relationships or activities that could appear to have influenced the submitted work. The views and opinions expressed in this article are those of the authors and do not necessarily reflect the official policy or position of the funder. The funding body has not had input into the design, data collection, analysis, interpretation of data, in the writing of the manuscript, nor in the decision to submit the manuscript for publication.

**Competing interests** PK, ACV, HE, KMMcP had financial support from the Ministry of Health of New Zealand for the submitted work; subsequent to the completion of this project and data analysis, KMMcP became the Chief Executive of the Health Research Council of New Zealand.

**Patient consent** Not required.

**Ethics approval** New Zealand Health and Disability Ethics Committee (Northern A, NTY/12/04/028/AM05) and the Auckland University of Technology's Ethics Committee (12/163).

**Provenance and peer review** Not commissioned; externally peer reviewed.

**Data sharing statement** Quantitative data from the study can be obtained from the author, subject to the funder's permission.

**Open Access** This is an Open Access article distributed in accordance with the Creative Commons Attribution Non Commercial (CC BY-NC 4.0) license, which permits others to distribute, remix, adapt, build upon this work non-commercially, and license their derivative works on different terms, provided the original work is properly cited and the use is non-commercial. See: <http://creativecommons.org/licenses/by-nc/4.0/>

© Article author(s) (or their employer(s) unless otherwise stated in the text of the article) 2018. All rights reserved. No commercial use is permitted unless otherwise expressly granted.

## REFERENCES

- Eivers AR, Brendgen M, Borge AIH. Stability and change in prosocial and antisocial behavior across the transition to school: teacher and peer perspectives. *Early Education & Development* 2010;21:843–64.
- White J, Connelly G, Thompson L, et al. Assessing wellbeing at school entry using the strengths and difficulties questionnaire: professional perspectives. *Educational Research* 2013;55:87–98.
- Kim-Cohen J, Caspi A, Moffitt TE, et al. Prior juvenile diagnoses in adults with mental disorder: developmental follow-back of a prospective-longitudinal cohort. *Arch Gen Psychiatry* 2003;60:709–17.
- Kim-Cohen J, Arseneault L, Newcombe R, et al. Five-year predictive validity of DSM-IV conduct disorder research diagnosis in 4(1/2)-5-year-old children. *Eur Child Adolesc Psychiatry* 2009;18:284–91.
- Bierman KL, Coie J, Dodge K, et al. School outcomes of aggressive-disruptive children: prediction from kindergarten risk factors and impact of the fast track prevention program. *Aggress Behav* 2013;39:114–30.
- Doughty C. *The effectiveness of mental health promotion, prevention and early intervention in children, adolescents and adults*. NZHTA Report. 8, 2005.
- Goodman R. The strengths and difficulties questionnaire: a research note. *J Child Psychol Psychiatry* 1997;38:581–6.
- Goodman R, Meltzer H, Bailey V. The strengths and difficulties questionnaire: a pilot study on the validity of the self-report version. *Eur Child Adolesc Psychiatry* 1998;7:125–30.
- Ministry of Health. *The B4 school check. A handbook for practitioners*. Wellington, 2008.
- Williamson A, D'Este C, Clapham K, et al. What are the factors associated with good mental health among Aboriginal children in urban New South Wales, Australia? Phase I findings from the Study of Environment on Aboriginal Resilience and Child Health (SEARCH). *BMJ Open* 2016;6:e011182.
- Embretson SE, Reise SP. *Item response theory for psychologists*. London: Lawrence Erlbaum Associates, Publishers, 2000.
- Cano S, Klassen AF, Scott A, et al. Health outcome and economic measurement in breast cancer surgery: challenges and opportunities. *Expert Rev Pharmacoecon Outcomes Res* 2010;10:583–94.
- Kersten P, Czuba K, McPherson K, et al. A systematic review of evidence for the psychometric properties of the Strengths and Difficulties Questionnaire. *Int J Behav Dev* 2016;40:64–75.
- Andrich D. A rating formulation for ordered response categories. *Psychometrika* 1978;43:561–73.
- Rasch G. *Probabilistic models for some intelligence and attainment tests (revised and expanded ed)*. Chicago: The University of Chicago Press, 1960/1980.
- Bond TG, Fox CM. *Applying the Rasch model. Fundamental measurement in the human sciences*. London: Lawrence Erlbaum Associates, 2001.
- Klein AM, Otto Y, Fuchs S, et al. Psychometric Properties of the Parent-Rated SDQ in Preschoolers. *European Journal of Psychological Assessment* 2013;29:96–104.
- Tobia V, Gabriele MA, Marzocchi GM. The Italian version of the Strengths and Difficulties Questionnaire (SDQ)-teacher: psychometric properties. *J Psychoeduc Assess* 2013;31:493–505.
- Mieloo CL, Bevaart F, Donker MC, et al. Validation of the SDQ in a multi-ethnic population of young children. *Eur J Public Health* 2014;24:26–32.
- Borg AM, Kaukonen P, Salmelin R, et al. Reliability of the strengths and difficulties questionnaire among Finnish 4–9-year-old children. *Nord J Psychiatry* 2012;66:403–13.
- Goodman A, Lamping DL, Ploubidis GB. When to use broader internalising and externalising subscales instead of the hypothesised five subscales on the Strengths and Difficulties Questionnaire (SDQ): data from British parents, teachers and children. *J Abnorm Child Psychol* 2010;38:1179–91.
- Wright BD. Comparing Rasch measurement and factor analysis. *Structural Equation Modeling: A Multidisciplinary Journal* 1996;3:3–24.
- Christensen KB, Engelhard J G, Salzberger T. Rasch vs. Factor Analysis. *Rasch Meas Trans* 2012;26:1373–8.
- Hagquist C. The psychometric properties of the self-reported SDQ – An analysis of Swedish data based on the Rasch model. *Pers Individ Dif* 2007;43:1289–301.
- Streiner DL, Norman GR. *Health Measurement Scales: a practical guide to their development and use*. Oxford University Press: Oxford, 2008.
- Statistics New Zealand. Māori population estimates: Mean year ended 31 december 2016. Secondary māori population estimates: Mean year ended 31 december 2016. 2016. [http://www.stats.govt.nz/browse\\_for\\_stats/population/estimates\\_and\\_projections/MaoriPopulationEstimates\\_HOTPMYe31Dec16.aspx](http://www.stats.govt.nz/browse_for_stats/population/estimates_and_projections/MaoriPopulationEstimates_HOTPMYe31Dec16.aspx)
- Høegh MC, Høegh SM. Trans-adapting outcome measures in rehabilitation: Cross-cultural issues. *Neuropsychol Rehabil* 2009;19:955–70.
- de Klerk G. Cross-Cultural Testing. In: Born M, Foxcroft CD, Butter R, eds. *Online readings in testing and assessment: international test commission*, 2008.
- Kersten P, Dudley M, Nayar S, et al. Cross-cultural acceptability and utility of the strengths and difficulties questionnaire: views of families. *BMC Psychiatry* 2016;16:347.
- Andrich D. *Rasch models for measurement series: quantitative applications in the social sciences no. 68*. London: Sage Publications, 1988.
- Linacre JM. Rasch power analysis: Size vs. Significance: Infit and outfit mean-square and standardized chi-square fit statistic. *Rasch Meas Trans* 2003;17:918.
- Tennant A, Pallant JF. The Root Mean Square Error of Approximation (RMSEA) as a supplementary statistic to determine fit to the Rasch model with large sample sizes. *Rasch Meas Trans* 2012;25:1348–9.
- Linacre JM. Sample size and item calibration [or Person Measure] stability. *Rasch Meas Trans* 1994;7:328.
- Wright BD, Tennant A. Sample size again. *Rasch Meas Trans* 1996;9:468.
- Charter RA, Feldt LS. Confidence intervals for true scores: is there a correct approach? *J Psychoeduc Assess* 2001;19:350–64.
- Nunnally JC, Bernstein IH. *Psychometric theory*. 3rd edn. New York: McGraw-Hill, 1994.
- RUMM2030 [program]. RUMM Laboratory Pty Ltd, 2009.
- Marais I, Andrich D. Effects of varying magnitude and patterns of response dependence in the unidimensional Rasch model. *J Appl Meas* 2008;9:105–24.
- Wainer H, Kiely GL. Item Clusters and Computerized Adaptive Testing: A Case for Testlets. *J Educ Meas* 1987;24:185–201.
- Grimby G. Useful reporting of DIF. *Rasch Measurement Transactions* 1998;12:651.
- Holland PW, Wainer H. *Differential Item Functioning*. NJ: Hillsdale: Lawrence Erlbaum, 1993.
- Smith EV. Detecting and evaluating the impact of multidimensionality using item fit statistics and principal component analysis of residuals. *J Appl Meas* 2002;3:205–31.
- Tennant A, Pallant JF. Unidimensionality matters! (a tale of two Smiths?). *Rasch Meas Trans* 2006;20:1048–51.
- Wright BD. Reliability and separation. *Rasch Meas Trans* 1996;9:472.
- Sheng Y, Sheng Z. Is coefficient alpha robust to non-normal data? *Front Psychol* 2012;3:1–13.
- Wright BD. Separation, reliability and skewed distributions: statistically different levels of performance. *Rasch Meas Trans* 2001;14:786.
- Masters GN. A rasch model for partial credit scoring. *Psychometrika* 1982;47:149–74.
- Stone LL, Otten R, Engels RC, et al. Psychometric properties of the parent and teacher versions of the strengths and difficulties questionnaire for 4- to 12-year-olds: a review. *Clin Child Fam Psychol Rev* 2010;13:254–74.
- Richter J, Sagatun Å, Heyerdahl S, et al. The Strengths and Difficulties Questionnaire (SDQ) - self-report. An analysis of its structure in a multiethnic urban adolescent sample. *J Child Psychol Psychiatry* 2011;52:1002–11.

50. Ortuño-Sierra J, Fonseca-Pedrero E, Aritio-Solana R, *et al.* New evidence of factor structure and measurement invariance of the SDQ across five European nations. *Eur Child Adolesc Psychiatry* 2015;24:1523–34.
51. Goodman A, Patel V, Leon DA. Why do British Indian children have an apparent mental health advantage? *Journal of Child Psychology and Psychiatry* 2010;51:1171–83.
52. Goodman A, Heiervang E, Fleitlich-Bilyk B, *et al.* Cross-national differences in questionnaires do not necessarily reflect comparable differences in disorder prevalence. *Soc Psychiatry Psychiatr Epidemiol* 2012;47:1321–31.
53. de Vries PJ, Davids EL, Mathews C, *et al.* Measuring adolescent mental health around the globe: psychometric properties of the self-report Strengths and Difficulties Questionnaire in South Africa, and comparison with UK, Australian and Chinese data. *Epidemiol Psychiatr Sci* 2017;30:1–12.
54. Kersten P, Vandal AC, Elder H, *et al.* Concurrent validity of the strengths and difficulties questionnaire in an indigenous pre-school population. *J Child Fam Stud* 2017;26:2126–35.
55. Williamson A, Redman S, Dadds M, *et al.* Acceptability of an emotional and behavioural screening tool for children in Aboriginal Community Controlled Health Services in urban NSW. *Aust N Z J Psychiatry* 2010;44:894–900.
56. Beaton DE, Bombardier C, Guillemin F, *et al.* Guidelines for the process of cross-cultural adaptation of self-report measures. *Spine* 2000;25:3186–91.
57. Samad L, Hollis C, Prince M, *et al.* Child and adolescent psychopathology in a developing country: testing the validity of the strengths and difficulties questionnaire (Urdu version). *Int J Methods Psychiatr Res* 2005;14:158–66.
58. Thabet AA, Stretch D, Vostanis P. Child mental health problems in Arab children: application of the strengths and difficulties questionnaire. *Int J Soc Psychiatry* 2000;46:266–80.
59. Lane C. *Adult literacy and numeracy in New Zealand – A regional analysis. Perspectives from the adult literacy and life skills survey, 2012.*
60. Statistics New Zealand. *2013 Census quickstats about culture and identity.* Wellington, New Zealand, 2014.
61. The Advisory Group on Conduct Problems. *Conduct problems. Best practice report.* Wellington, 2009.
62. Statistics New Zealand. Te Kupenga. 2013.
63. Hedley C, Thompson S, Morris Mathews K, *et al.* The B4 school check behaviour measures: findings from the hawke's bay evaluation. *Nursing Praxis in New Zealand* 2012;28:13–23.