

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21

Developing a Rasch Measurement Physical Fitness Scale
for Hong Kong Primary School-aged Students

Zi Yan

The Hong Kong Institute of Education

Trevor G. Bond

James Cook University

Request for reprints should be addressed to

Dr. Zi Yan

Address: Centre for Special Needs and Studies in Inclusive Education, The Hong Kong Institute of
Education, 10 Lo Ping Road, Tai Po, N.T., Hong Kong.

Phone: (+852) 2948 6367

Email: zyan@ied.edu.hk

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Abstract

The main purpose of this study was to develop a Rasch Measurement Physical Fitness Scale (RMPFS) based on physical fitness indicators routinely used in Hong Kong primary schools.

A total of 9,439 records of students' performances on physical fitness indicators, retrieved from the database of a Hong Kong primary school, were used to develop the Rasch scale.

Following a series of iterative Rasch analyses which adopted the "data should fit the model"

approach, four physical fitness indicators (i.e., 6-minute Run, 9-minute Run, 1-minute

Sit-ups and Dominant Handgrip) were successfully calibrated to form the RMPFS. The

RMPFS and its scale indicators showed fit to the Rasch model sufficient for the intended

purposes of measuring overall fitness of children. The overall physical fitness measure

reflects children's fitness on three key core components of physical fitness (i.e.,

cardiorespiratory fitness, muscular endurance, and muscular strength). Advantages of the

RMPFS are discussed and recommendations for future research follow. The findings of this

study provide a better knowledge basis for interpreting children's physical fitness assessment

results.

Key words: Rasch measurement, physical fitness, primary school, data should fit the model

1 Developing a Rasch Measurement Physical Fitness Scale

2 for Hong Kong Primary School-aged Students

4 **Introduction**

5 Given the important role physical fitness should play in children's lives, fitness
6 assessment/testing is intuitively a crucial part of physical education which aims to promote a healthy and
7 physically active lifestyle. However, fitness testing in schools has being criticised over decades, and even
8 its necessity for children has been seriously questioned (Liu, 2008). The special issue on youth fitness
9 testing published in *Measurement in Physical Education and Exercise Science* (MPEES) in 2008
10 thoroughly discussed different perspectives on youth fitness testing. For example, a pedagogical
11 perspective argued that fitness tests should be implemented as formative evaluation. Then fitness testing
12 results should be informative for teaching and learning in physical education (Silverman, Keating, &
13 Phillips, 2008). In terms of promoting physical activity, fitness assessments are expected to provide
14 accurate measures carrying important information about children's health-related fitness levels. Therefore,
15 they could optimize the effectiveness of physical education (Welk, 2008). Moreover, there is no doubt
16 that the use and interpretation of fitness assessment have important educational, pedagogical, and
17 psychological consequences (Mahar & Rowe, 2008). In summary, the editors and authors of the MPEES
18 special issue agreed that youth fitness testing can serve a useful purpose in school settings if used in the
19 correct way. This article aims to extend this "correct way" discussion by shedding some light on how to
20 achieve objective physical fitness measurement based on fitness testing scores.

1 Accurate measures of youth fitness are needed by both researchers and educators regardless of
2 their purposes (Mahar & Rowe, 2008). The routine practice in traditional approaches is that different
3 components of physical fitness (e.g., body composition, cardiorespiratory fitness, flexibility, muscular
4 endurance, and muscular strength) are assessed using different indicators and children's abilities in each
5 of these components are reported and interpreted using raw scores (in metres, kilograms, seconds, etc.) or
6 percentile ranks. However, raw scores might not provide a valid *measure* because they have little
7 inferential value (Wright, 1997; Wright & Mok, 2000). The validity of raw scores in representing fitness
8 levels in this approach is based on an unquestioned assumption, namely, the raw scores are accepted
9 implicitly as being equal interval. Unfortunately, the raw scores themselves (unless used to derive further
10 criterion measures, e.g., estimated VO₂max based on scores in 6/9-minute Run test) actually indicate only
11 the ordering of the children's performances, but have little inferential value about the size of the
12 differences among scores in terms of "fitness". While metres indicate equal amounts of difference on the
13 length or distance scales, it is an act of faith to conclude that metres indicate equal difference on the
14 cardiorespiratory fitness scale. Metres have only ordinal meaning when they are used as the score units in
15 the 6-minute Run test, therefore they might not yield valid measures of the underlying fitness component.

16 Another deficiency associated with the traditional approaches to physical fitness assessment is
17 that the interpretation of results of physical fitness assessment in norm-referenced framework is often not
18 accurate nor comparable because of the sample-dependence and indicator-dependence of assessments,
19 where ranks or percentiles are provided in interpreting students' performances on physical fitness
20 indicators. Those ranks or percentiles provide only an inexact basis for comparison among students and,
21 rather, should be regarded as indicators of students' relative strengths and weaknesses (Williams,
22 Harageones, Johnson, & Smith, 2000). However, use of raw numbers/counts and the allocation of

1 norm-referenced ratings do not allow for the direct assessment of children's fitness against some objective
2 fitness standard in which the measurement and interpretation of students' fitness levels is independent of
3 sample and indicator.

4 Furthermore, it is time-consuming using the traditional approach to administering all fitness tests
5 to the whole class with 40 or more students. Since a single total score might not provide a meaningful
6 summary of different fitness indicators, multifaceted profiles which contain scores for each component of
7 physical fitness are often regarded as more appropriate (Fleishman, 1964; Marsh 1993). A consequent
8 by-product is that assessment tasks in the physical education curriculum increase teachers' workload and
9 occupy resources which could be put into teaching. There is little doubt that physical fitness is a
10 multifaceted concept, but the extent to which any set of multidimensional indices used in traditional
11 approaches should disqualify a unidimensional fitness index still remains open for discussion as well as
12 evidence-based empirical investigation. The question addressed in this article is to what extent is it
13 possible to generate a unidimensional index of physical fitness, which provides interval scale fitness
14 measures for children independent of sample and indicator, for estimating differences between groups of
15 children and for tracking changes in fitness levels over time.

16 The Rasch model (Andrich, 1988; Rasch, 1960) provides ways to address the deficiencies
17 inherent in traditional approaches to physical fitness assessment. Firstly, Rasch analysis can transform
18 non-linear raw scores into logit scale measures which have constant interval meaning and provide
19 objective and linear measurement from ordered category responses (Linacre, 2000, 2006a). Secondly, the
20 feature of "parameter separation" or "invariance of parameters" (Bond & Fox, 2007; Wright & Masters,
21 1982) of the Rasch model implies that the calibration of fitness indicators is sample-distribution free and
22 the calibration of persons is indicator-distribution free along the fitness continuum. The

1 sample-distribution free calibration of fitness indicators means that the difficulty estimates of indicators
2 (e.g., 6-minute Run, 1-minute Sit-ups, etc.) should be invariant, within measurement error, no matter
3 which sample is used to calibrate those indicators. The indicator-distribution free calibration of persons
4 means that the fitness estimate of any person should remain invariant, within measurement error, no
5 matter which particular fitness indicators are used to measure that person's fitness. Therefore, direct
6 person-person, item-item, and person-item comparisons can be conducted easily, based on their locations
7 on the common logit scale. Finally, an overall fitness measure can be provided for a student, even if s/he
8 had not performed on all of the physical fitness indicators which have been calibrated onto the fitness trait
9 continuum.

10 Unlike more general multidimensional or IRT models and other (true score) statistical techniques
11 that adopt a "the model fits the data" approach and manipulate the different parameters to accommodate
12 the idiosyncrasies of any data set, the Rasch model requires that "data fit the model" (Andrich, 2004) for
13 the purpose of achieving objective measurement. This is one of the key differences between Rasch-based
14 studies and other quantitative studies in the human sciences. The Rasch model is held as being able to
15 solve the basic measurement problem common to all social sciences (Andersen, 1995) and it has been
16 applied in sport sciences and physical education studies by a growing number of researchers whose
17 reviews provide more detail (e.g., Strauss, Büsch, & Tenenbaum, 2007; Tenenbaum, Strauss, & Büsch,
18 2007). For example, Rasch analysis has been utilized to calibrate physical function or competence (Zhu &
19 Kurz, 1994), perception of sports games (Kang & Kang, 2006), and difficulty levels of physical fitness
20 indicators (Zhu & Safrit, 1993). Studies have applied the Rasch model to develop or evaluate instruments
21 used in exercise studies. Hands & Larkin (2001) studied children's performance on different motor tasks
22 and developed two separate unidimensional Rasch scales of motor abilities for boys and girls respectively.

1 Zhu, Timm and Ainsworth (2001) modified an exercise barriers instrument and validated it using the
2 Rasch model framework. Heesch, Masse and Dunn (2006) used Rasch analysis to re-evaluate three
3 commonly used scales including the Physical Activity Enjoyment Scale, the Benefits of Physical Activity
4 Scale and the Barriers to Physical Activity Scale. B üsch et al. (2009) used a mixed Rasch model to
5 investigate the construct validity of the German general motor fitness and coordination test for children.
6 They found that two qualitatively different classes of children could be distinguished. Members of the first
7 class were characterized by high running ability and low throwing ability whereas members of the second
8 class were characterized by low running and high throwing abilities.

9 Tenenbaum, Strauss, and B üsch (2007) claim that the application of Rasch model in physical
10 education and sport sciences is promising from both a methodological and a content-related perspective.
11 A number of advantages of Rasch model analyses have been echoed in previous studies. In calibrating a
12 gross motor skills instrument with the many-facets Rasch model, Zhu and Cole (1996) demonstrated the
13 advantages of the Rasch model over the traditional norm-referenced interpretation, including benefits of
14 parameter separation, sharing the same metric among items and examinees, and providing linear measures.
15 They also pointed out that the person measures, together with S.E. and fit statistics, provided useful
16 diagnostic information to identify strengths and weaknesses of examinees. Bowles and Ram (2006)
17 revealed that Rasch analyses of volleyball players' performances on three skills (serve, serve receive, and
18 attack) produced an equal-interval scale which provided more objective and consistent information about
19 volleyball players' abilities than could be obtained by traditional instruments. Zhu (2001) found that
20 Rasch model could accurately equate different motor function tests so that cross-test scores could be
21 interpreted in a common measurement framework, an important outcome which remains unachievable in
22 traditional approaches to motor function assessment. B üsch and Strauss (2005) used the Rasch

1 measurement model to study 503 participants' performance on 6 gross-motor coordination tasks,
2 categorized as precision and time-pressure tasks. They found that persons performing gross-motor
3 coordination tasks could be differentiated based on the coordination strategy they used. The results
4 displayed the advantages of Rasch model in identifying strategies used by persons in completing
5 gross-motor tasks and distinguishing between person and item characteristics.

6 The Rasch model has also been applied in attempts to combine closely related scales to assess
7 single unidimensional physical functioning constructs. An interesting study conducted in the health care
8 domain combined two separate but related scales into one unidimensional scale (Hsueh, Wang, Sheu, &
9 Hsieh, 2004). The 10-item Barthel Index (BI) assessing Activities of Daily Living (ADL) and the 15-item
10 Frenchay Activities Index (FAI) assessing instrumental ADL were administered to 245 patients one year
11 after stroke. The data from these two scales were combined and analyzed using the Rasch model and the
12 result indicated that all but two FAI items fit the unidimensional Rasch model very well, indicating that
13 the BI and the FAI assess a single underlying unidimensional ADL construct. Further analyses of the
14 23-item unidimensional scale revealed that it had quite high person reliability (0.94) and the range of item
15 difficulties was well targeted to the patient sample. A "look-up" conversion table was then offered to
16 transform combined BI and FAI raw scores into Rasch ADL interval measures. Thus a clinically useful
17 instrument was developed by combining the BI and the FAI scales and the new scale had improved range
18 and sensitivity for assessing comprehensive ADL function.

19 However, this kind of combining attempt is seldom found in physical education literature.
20 Traditional approaches conceptualize physical fitness as a multifaceted construct, hence psychometrically
21 multidimensional. However, a single *overall* fitness score is still preferable, even necessary in many
22 situations, especially for the interpretation of students' comprehensive physical ability. In most cases, an

1 overall fitness score is obtained by simply summing or averaging the scores for different components of
2 physical fitness. It is obvious that such averaged overall fitness scores lose quite large amounts of
3 information about specific fitness aspects and, therefore, one should be very cautious in interpreting that
4 kind of overall score (Fleishman, 1964). Furthermore, as argued by B üsch et al. (2009), item homogeneity
5 needs to be checked before using a sum score to estimate general fitness. Marsh (1993) recommended
6 constructing, if necessary, a weighted summary score that assigns an optimal weight to each component
7 based on theoretical and empirical research. But it remains a considerable challenge to derive optimal
8 weights for different fitness components because the weights might need modification according to
9 particular criteria, particular research purposes, or even the predisposition of the particular investigator.

10 Thus, the main purpose of this study was to develop, to the extent that it was both useful and
11 possible, a Rasch Measurement Physical Fitness Scale (RMPFS) combining all, or at least some of the
12 indicators routinely used in Hong Kong primary schools. A successful scale would then calibrate person
13 ability (students' overall physical fitness levels) and item difficulty (difficulty levels of each of the
14 physical fitness indicators) in a single, stable fitness measurement framework. Given the review of the
15 quantitative approaches open for adoption in such a research project, the position taken in this research is
16 the primacy of the requirement to produce scientifically repeatable measures based on the principles
17 espoused in Rasch measurement. As a consequence, this particular research explicitly adopts the Rasch
18 'data fit the model' approach for the empirical investigation of the construction of physical fitness
19 measures for children.

20

1 **Method**

2 **Data**

3 Fitness data used in this study were retrieved from the physical fitness assessment records
4 database of a large, regional Hong Kong primary school, a government-subsidized primary school located
5 in the north-eastern New Territories of Hong Kong. This school routinely has five classes at each year
6 level from primary 1 (6 years old) to primary 6 (12 years old) with an annual enrolment of over 1,000
7 students.

8 The data set covers this school's students' physical fitness records for the academic years
9 2002-03 to 2006-07. There are two rounds of students' records for each academic year, except 2002-03
10 for which the records for the 2nd semester were not entered into the school's database. Initially, 10,512
11 student records were included in the potential data pool for this study, and finally 9,439 records were kept
12 for scale development after excluding exceptional and unreasonably extreme data. It is worth pointing out
13 that each record does not necessarily refer to an independent student since this is a longitudinal data set
14 over five years and most students would have several records (potentially up to nine) over time in the data
15 set. Of the records, there are 5,149 records (54.6%) for males and 4,290 records (45.4%) for females. The
16 age range for all records extends from 6 to 13 years ($M = 8.53$, $SD = 1.73$). Only four records did not
17 include age information. The details of the sample used in this study are presented in Table 1.

18

19

Put Table 1 about here

1 **Physical Fitness Indicators**

2 The partner school of this study administers the physical fitness testing recommended by the
3 School Physical Fitness Award Scheme (Hong Kong Education and Manpower Bureau, 2005) except that
4 body composition is estimated by BMI and not the skinfold method, because the equipment required for
5 skinfold testing is not available in this school. The other eight fitness indicators include 6-minute Run,
6 9-minute Run, 1-minute Sit-ups, Standard Push-ups, Modified Push-ups, Right Handgrip, Left Handgrip,
7 and Sit-and-Reach. It is worth noting that the 6-minute Run test is administered to grades 1 to 3 students
8 only and the 9-minute Run test is administered to grades 4 to 6 students only. The Standard Push-ups test
9 is administered to grades 3 to 6 male students only and the Modified Push-ups test is administered to all
10 grades 1 and 2 students as well as grades 3 to 6 female students.

11 **Data Analysis**

12 The software package used for Rasch analyses in the present study is WINSTEPS 3.0 programme
13 (Linacre, 2006a) and the Partial Credit Model (PCM) was specified for these analyses¹. The PCM is a
14 sound option for Rasch analyses with the physical fitness data in this study considering that the definition
15 of the rating scale is unique for each of the physical fitness indicators. The “partially correct response(s)”
16 between incorrect and completely correct item responses provided in the PCM are in accordance with the
17 different levels of performances between minimum and maximum raw scores on each physical fitness
18 test.

¹ Although the Rasch Poisson Counts Model has been used to measure physical fitness (e.g., Zhu & Safrit, 1993), the appropriateness of the Rasch Poisson Counts Model in time-limited psychomotor performances such as 1-minute Sit-ups test scores is dubious, because some of the model’s requirements are not satisfied by such data (Zhu & Safrit, 1993). For example, the Rasch Poisson Counts Model assumes that examinees should complete the repetitions at a constant speed through the whole performance. However, the effect caused by fatigue in the 1-minute Sit-ups test violates this basic assumption; repetition speed is usually slower and slower as examinees complete greater numbers of sit-ups during the one minute period.

1 **Logarithmic Transformation of Raw Scores**

2 Although it is not a problem for WINSTEPS to handle data with a large number of ordered
3 response categories (it accommodates up to 255 category levels per item), using many more than the
4 necessary category levels is likely to introduce challenges to the meaningful interpretation of the results.
5 From a practical perspective, it is unlikely that primary school-aged students' performances on physical
6 fitness indicators have more than about 10 useful qualitatively different levels: it is unlikely that 10
7 metres in a 6-minute Run test or one centimetre in a Sit-and-Reach test indicate meaningful differences in
8 overall physical fitness levels, even if such a small difference could move a child's fitness estimate from a
9 lower to a higher response category for that one indicator. Thus re-expressing raw data into a reasonable
10 number of ordered categories would help the interpretation of the results and the detection of departures
11 from fit to the model more clearly (Linacre, 2000). A Poisson logarithmic transformation was used to
12 transform the raw scores into a data set with more even distribution and more meaningful category
13 structure. The transformation can be expressed as

$$14 \quad \text{Scored category} = 1 + 8 * \frac{\log(\text{observation}+1) - \log(L+1)}{\log(H+1) - \log(L+1)} \quad (1)$$

16 Where L is the lowest value of the observations, and H is the highest value of the observations. The
17 number "8" was chosen just because after some initial investigation analyses a 9-category structure was
18 selected as the appropriate transformation target.

19 **Iterative Sequence of Analytical Steps**

20 Given that the RMPFS would be developed from a Rasch measurement perspective, each fitness
21 indicator where data violated Rasch measurement requirements was excluded in turn from the scale. More

- 1 specifically, this study took the strong “data fit the model” approach in developing this physical fitness
2 scale. Seven criteria were utilized in the procedure of scale development to investigate the quality of those
3 indicators and to decide whether an indicator should be retained in or excluded from subsequent analyses.
- 4 ● *Investigations from a practical perspective.* Practical considerations undertaken before
5 undertaking statistical analyses might uncover some important factors detrimental to scale
6 development.
 - 7 ● *Fit statistics for indicators.* Since Mean Square (MNSQ) statistics are relatively more stable than
8 are their standardized forms (ZSTD) in Rating Scale Model and Partial Credit Model analyses
9 (Smith, Rush, Fallowfield, Velikova, & Sharpe, 2008), MNSQs are used as fit criteria for scale
10 quality assurance.
 - 11 ● *Point-measure correlations for indicators.* The point-measure correlation coefficient of an
12 indicator refers to the correlation between the estimate for any particular fitness indicator and the
13 overall measure of the fitness trait under measurement (Linacre, 2006b). Normally, a
14 point-measure correlation coefficient higher than +0.4 indicates acceptable consistency of
15 indicator polarity in a scale.
 - 16 ● *Rasch reliability.* Rasch measurement provides both person reliability and item reliability indices.
17 Rasch person reliability refers to the consistency of person ordering along the trait continuum
18 measured by the scale (Smith, 2001; Wright & Masters, 1982), and Rasch item reliability
19 indicates replicability of item placements along the trait continuum if the same set of items were
20 administered to another similar sample of persons (Bond & Fox, 2007).
 - 21 ● *Variance explained by measures.* Variance explained by the measures refers to the proportion of
22 variance of observations that could be explained by the item difficulties, person abilities and

1 rating scale structures in Rasch analyses (Linacre, 2006a). A higher proportion of variance
2 explained by Rasch measures means that the Rasch model has better capacity for predicting
3 performances of both items and persons.

4 ● *Response category structure.* Successful implementation of polytomous Rasch measurement
5 requires well functioning performance categories for each indicator in the scale.

6 ● *Influence of underfitting persons.* The impact of extremely misfitting persons, especially
7 underfitting (erratically performing) persons, on fitness scale quality would be investigated.

8

Results

Through a theory-driven iterative developmental procedure guided by Rasch model measurement perspectives as well as practical considerations, eventually four physical fitness indicators including 6-minute Run, 9-minute Run, 1-minute Sit-ups and Dominant (not Left- or Right-) Handgrip were successfully calibrated to form the RMPFS; thereby integrating three key components of physical fitness - cardiorespiratory fitness, muscular endurance, and muscular strength - to provide a single person measure of overall physical fitness suitable for use with primary school children in Hong Kong.

The indicator properties of the RMPFS are presented in Table 2. The difficulty levels (i.e. item measures) for the four fitness indicators range from -1.59 logits to +1.25 logits associated with standard errors of 0.02 or 0.03 logits. These small standard errors imply that the indicator difficulty estimations are quite precise – primarily due to the large calibration sample. The Infit and Outfit MNSQs range from 0.85 to 1.13, indicating sufficient fit to the Rasch model for practical measurement purposes – especial for such low-stakes decisions as monitoring children’s fitness levels in school settings. The point-measure correlations approximate 0.8, supporting the claim that all the indicators function in the same direction as a part of the physical fitness latent trait under measurement. The Rasch item reliability is 1.00 and the Rasch person reliability is lower, but acceptable, at 0.77, a consequence of retaining only four indicators in the RMPFS.

Put Table 2 about here

Figure 1 presents the Wright map of the 4-indicator RMPFS. Students are placed on the left side of the scale according to physical fitness, and the fitness indicators are shown on the right side. The

1 students with the highest fitness levels and the fitness indicators with highest difficulty levels are located
2 at the top of the map, while the students with the lowest fitness level and the easiest fitness indicators are
3 located at the bottom. The means of student measures and indicator difficulty calibrations are shown as
4 the corresponding *M*s on the map. The *S*s and *T*s represent ± 1 and ± 2 standard deviations of the student
5 and indicator distributions respectively. It can be seen that the difficulty levels of the RMPFS physical
6 fitness indicators ($M = 0.00$, $SD = 1.16$) are appropriate for these students' fitness levels ($M = -0.21$, $SD =$
7 2.78). The range of indicators' difficulty (-1.59 to 1.25 logits) is much smaller than the range of students'
8 ability (-12.86 to 11.17 logits). However, the ranges of difficulty levels of the response categories
9 (categories 1 to 7) for each indicator, as presented on the right-hand side of the map (from R6_1 at -11.31
10 logits to DH_7 at $+10.94$ logits), reveals that the indicators overall provide good coverage of the fitness of
11 the primary school-aged students in this sample.

13 Put Figure 1 about here

15 The Item Characteristic Curves (ICCs) and category probability curves provide further support
16 for the valid functioning of the scale. Figure 2 presents the empirical and expected ICCs for the four
17 indicators. It can be seen that the empirical ICCs match the theoretical ICCs reasonably well, especially
18 for students' with median fitness levels located around the middle of the curves. There are larger
19 discrepancies between the empirical and theoretical ICCs for the most able and the least able students
20 located at the extremes of the curves.

22 Put Figure 2 about here

1 The category probability curves for each of the four indicators presented in Figure 3 show that
2 each performance category has a distinct peak in the graph for all indicators. That means each category
3 for each indicator was the most probable performance level for given groups of persons with a specific
4 level of physical fitness. There is no evidence of category threshold disordering (Bond & Fox, 2007;
5 Linacre, 2002) and the threshold calibrations advance monotonically with category, indicating that higher
6 performance categories correspond to higher measures of physical fitness.

7

8

Put Figure 3 about here

9

10

Discussion

Initially a total of nine physical fitness indicators were used to develop the RMPFS, but only four of them were retained to form the final interval-level measurement scale. The other five indicators were excluded or replaced in the development procedure based on the seven criteria described in the method section. The development procedure is presented in Table 3. Scales 1 to 5 displayed in Table 3 are intermediate scales before Scale 6 - the final version of RMPFS – was finally established.

Put Table 3 about here

Consideration of BMI

BMI was excluded for conceptual and practical reasons. The first, BMI is a rough index appropriate for reporting adiposity at the population level but not optimal for use with individuals because of prediction error (Heyward, 2002; Stratton & Williams, 2007). The second, BMI is a trait with an inverted U-shaped (\cap) distribution. A higher BMI score does not necessarily stand for a better level of physical fitness; nor does a lower BMI score. This is a distinctive feature which sets BMI apart from other fitness indicators. Combining BMI together with other indicators in the Rasch measurement scale would contradict one of the requirements of Rasch model: all items in the same scale should function in the same (linear) direction along the latent trait under measure.

Consideration of Sit-and-Reach

Sit-and-Reach, which is used to assess flexibility, is distinct from the other indicators in some important ways. Students' performances for other indicators increase monotonically with students' age, but it is not the case for Sit-and-Reach. Furthermore, the correlation matrix shows that flexibility

1 component has relatively low correlations with other indicators of physical fitness. This is consistent with
2 the findings of other studies: Marsh and Redmayne (1994) found that the correlations involving the
3 flexibility component are smaller than the correlations involving other components of physical fitness.
4 Therefore, Sit-and-Reach was excluded from the RMPFS to see if there was any subsequent improvement
5 in scale properties. The results in Table 3 show that the Rasch person reliability increased appreciably
6 from 0.52 to 0.66 even though the raw score range of the scale was reduced.

7 **Consideration of Handgrip**

8 Rasch factor analyses of fit residuals show that Right Handgrip and Left Handgrip have quite high
9 loadings on the 1st contrast factor and the correlation between their residuals is 0.52, i.e., they share about
10 27% of their variance in common. That suggests there is probably a separate fitness sub-dimension
11 comprising of Right Handgrip and Left Handgrip and that there is local dependency between these two
12 indicators. From the Rasch perspective, one promising solution is to use Dominant Handgrip instead of
13 Right Handgrip and Left Handgrip. In this case, the higher score of Right Handgrip and Left Handgrip
14 was chosen as the Dominant Handgrip result for each student. As well as for Right Handgrip and Left
15 Handgrip, local dependence is likely to occur between 6-minute Run and 9-minute Run or between
16 Standard Push-ups and Modified Push-ups considering their very similar nature. However, this is not a
17 concern for the current analyses since no single case in the data set has scores on both 6-minute Run and
18 9-minute Run or Standard Push-ups and Modified Push-ups.

19 **Consideration of Push-ups**

20 The properties of Scale 3 (see Table 3) show that the Standard Push-ups and Modified Push-ups
21 have poor fit to the Rasch model. There are two reasons that probably introduce noise to these two

1 indicators. The first, these two indicators are usually used for secondary school-aged students in Hong
2 Kong, but not for primary school-aged students. The partner school of this study used these two indicators
3 just as supplementary tests for a small portion of students (14.4% for Standard Push-ups and 20.2% for
4 Modified Push-ups) before academic year 2005-06. The second reason is related to the nature of the
5 Push-ups test. These two tests have no time limit but have an assumption about students' willingness to
6 participate, i.e., students were assumed to try their best to complete as many push-ups as possible (until
7 they cannot do any more). But this not always seem to be the case in practice, especially for
8 supplementary tests to which students often attach less importance. Considering the misfit shown by these
9 two indicators and the possibility of measurement noise introduced by them, it is reasonable to exclude
10 them from further RMPFS development. As indicated in Table 3, the properties of 4-indicator scale
11 (Scale 4) are much better than those of previous versions. The Rasch person reliability increased from
12 0.60 to 0.63. The variance explained by measures increased considerably from 62.6% to 66.9%.

13 **Optimizing Response Category**

14 The results of Rasch analyses adopting 9-category data set showed that the response category
15 structure was not optimal because 1) the distribution of respondents among categories was not even; 2)
16 there were some reversed average measures and threshold calibrations; and 3) the category probability
17 curves for some categories were submerged by others. Therefore, it is appropriate to collapse some
18 adjacent and potentially redundant categories in order to obtain a meaningful and interpretable category
19 structure for each indicator. Two principles were followed in the process of combining adjacent
20 categories. The first was to ensure each category had a reasonable number of respondents, and the second
21 was to ensure average measures for categories and threshold difficulties increase monotonically and with

1 reasonable increments. Finally, a 7-category structure was developed and the category functioning
2 effectiveness was examined in detail according to the guidelines suggested by Linacre (2002). The
3 point-measure correlation coefficients of all the 4 indicators range from 0.71 to 0.79. The number of
4 observations of all categories for each indicator ranges between 14 and 3,496 with a mean of 879. The
5 observation distributions across categories for all indicators are unimodal distributions peaking in a
6 central category and show smooth decreases to category 1 and category 7 respectively. The average
7 measures of categories for all indicators advance monotonically with category. The Outfit and Infit
8 MNSQs for all categories range between 0.79 and 1.44, and most of them are very close to 1.0. The
9 threshold calibrations of categories for all indicators advance monotonically. The measure-to-category
10 coherence and category-to-measure coherence for most of the categories except categories 1 to 7 are
11 acceptable. The distances between adjacent threshold calibrations are all larger than 1.0 logit and less than
12 5.0 logits with only two exceptions. Therefore, Scale 5 based on 7-category data replaced Scale 4 which
13 had been constructed from 9-category data (see Table 3).

14 **Influence of Underfitting Persons on the RMPFS**

15 Verhelst and Glas (1995) stated that there are two methods to improve Rasch measurement scale
16 construction. The one is to eliminate “bad” items; the other is to exclude temporarily some test takers
17 whose performances do not fit the Rasch model. At this point, eliminating further items from the scale is
18 not the preferable option because only four indicators are retained in Scale 5 (see Table 3) and all of them
19 are “acceptably good” items from both practical and Rasch measurement perspectives. Consequently, the
20 alternative – temporarily eliminating misfitting persons who introduce unexpected noise to the
21 measurement – was carried out in order to investigate possible improvement in the measurement

1 characteristics of the scale. Bond and Fox (2007) pointed out that underfitting persons ($MNSQ \gg 1.0$)
2 are more detrimental to calibrating a measurement scale than are overfitting persons ($MNSQ \ll 1.0$).
3 Linacre (2002) further stated that MNSQs higher than 2.0 indicate more noise than useful information
4 provided by the observations. Consequently, persons were excluded from the scale construction if either
5 their Outfit MNSQ or Infit MNSQ was higher than 2.0 on Scale 5. Finally, a total of 1,185 cases were
6 excluded and the final version of RMPFS was established based on the retained 8,469 cases which had at
7 least one score for any of the four indicators (6-minute Run, 9-minute Run, 1-minute Sit-ups, and
8 Dominant Handgrip). The results in Table 3 showed that the scale constructed without underfitting
9 persons exhibits significant improvement in both Rasch person reliability (increased from 0.62 to 0.77)
10 and variance explained by measures (increased from 68.7% to 81.5%).

11 **Properties of the RMPFS with Subsamples**

12 As described before, the data used to develop the RMPFS came from a longitudinal data set
13 collected over five years and most students would have several records over time in the data set. That
14 means some records might be considered as dependent on each other as they are the performances of the
15 same student at different time points. The reason for including all data in the primary calibration analysis
16 was to develop the RMPFS with as much good quality data as possible. At this point the concerned reader
17 might have some reservations due to the nature of the complete sample. Rasch modeling requires
18 performances from persons who are independent of each other. Otherwise the attractive feature of sample
19 distribution free measurement as well as the property of local independence of Rasch models might be
20 lost. In the complete analysis each record has been treated as that of an individual person. To rule out any
21 concern in that regard, we have completed separate Rasch analyses for the 4 items test (RMPFS) using

1 subsamples in which each person has only one record (e.g., potentially for the 9 separate subsamples for
2 each measurement point in time (approx. 1,000 pupils each). This approach examines the robustness of
3 the RMPFS to these apparently dependent records. Eventually, six subsamples with records for all four
4 RMPFS indicators were used. The subsamples for 1st semester of academic year 2002-03, 2nd semester
5 of academic year 2003-04, and 2nd semester of academic year 2004-05 were excluded since the records
6 for one or more of RMPFS indicators were missing. The results are presented in Table 4.

7
8 Put Table 4 about here

9
10 It can be seen from Table 4 that the properties of the RMPFS with each of the six independent
11 subsamples are quite good. The Infit and Outfit MNSQs range from 0.74 to 1.19. The point-measure
12 correlations range from 0.69 to 0.87. The Rasch item reliability is 1.00 and the Rasch person reliability
13 approximates 0.8. The standard errors of item estimates are, although still quite small, slightly larger than
14 those derived from the overall data due to the decrease of the sample sizes. The ordering of the item
15 difficulty is consistent with that for the overall sample with only one exception (2003-1: the subsample
16 from the first semester of academic year 2003-04) for which Dominant Handgrip appeared to be easier
17 than 6-minute Run, whereas that is not the case for other subsamples. Given that the item ordering
18 remained invariant except in this one instance, the concern about potential lack of person independence in
19 the results of the analysis for the whole data set can be put to one side.

20 **Age-dependent or Age-related?**

21 At this point, it could be easy to conclude that the RMPFS is merely reflecting changes in
22 children's body and fitness that are determined by their age. Figure 4a reflects the differences in fitness

1 levels, on average, between boys and girls at each age levels from 6 to 11 years age. While the differences
2 and trends are relatively consistent, the overlap across sexes and age groups remains quite substantial.
3 Close inspection of Figures 4b and 4c reveal the subtle changes that are not easily discerned in the
4 summary tables. The inference to be drawn from these graphs is, rather, that increases in children's fitness
5 is merely related to, but not actually determined by or dependent on children's age. If it were the case that
6 physical fitness in children was age-dependent, rather than age-related, it should be possible to determine
7 any children's age by his/her location on the fitness scale, or his/her fitness score simply referring to
8 his/her age. However, figures showing the full distribution of fitness scores by ages for boys and girls
9 reveal that this is clearly not the case. A boy with an RMPFS measure of 3.17 logits might be an 11-year
10 old of above average fitness for his age, a 9-year old who is a little fitter than average or, indeed, the
11 fittest boy in the first year of primary school, aged 6 years. Further, from Figure 4c, a 9-year old girl
12 (except for the very fittest of that age group) could have a fitness level that appears in the plots at any age
13 group of girls from 6-years old to 11-years old.

14
15 Put Figure 4 about here

16
17 It is obvious that students of Grades 1 to 3 (arguably less fit) will not be administered the more
18 difficult 9-minute Run test, while students of Grades 4 to 6 (apparently more fit) will not be administered
19 the less demanding 6-minute Run test. But because the scale category calibrations of the R6 and R9 are
20 quite close to each other (as shown in Figure 1), it might give the impression that by holding SU and DH
21 performances constant, those who scored a 4 on the R6 scale will have the same (or more or less the same)
22 abilities as those who scored a 4 on the R9 scale. Is it then reasonable to conclude that a 2nd grader who

1 took the 6-minute Run test will perform similarly (i.e. have similar level of physical fitness) to a 5th
2 grader who took the 9-minute Run test if that more difficult test was administered? Conversely, would a
3 5th grader perform similarly to a 2nd grader in the 6-minute Run test if given the easier test? The *prima*
4 *facie* evidence to support the category equivalence conclusion for R6 and R9 is displayed as Table 5:
5 although the *number* of meters covered varies by category according to whether R6 or R9 was
6 administered, the *speed* (meters per min) varies by category, but is independent of actual test. Given that
7 the two sub-samples, Grades 1 to 3 and 4 to 6, were not independently calibrated and then equated, it
8 remains open to future investigation concerning the extent to which their person measures will be useful
9 for tracking changes in fitness levels over time.

10

11

Put Table 5 about here

12

13

14

1 fitness measure that summarizes a student's physical fitness in different components. Even if the student
2 had not performed on all four RMPFS physical fitness indicators, that student can be given an overall
3 RMPFS fitness measure. By constructing an objective measurement scale for overall fitness, we can
4 locate the fitness of the students on an interval level measurement scale which of course maintains the
5 lower level ordinal relationships between the positions of all students. This means that any fitness level
6 (e.g., +1.0 logit) is objective, i.e., independent of any personal characteristic (e.g., sex or age) of the child.
7 Every child with RMPFS measure of +1.0 logit has the same fitness level (within error) which is 2 logits
8 higher than any child with a -1.0 logit RMPFS score and 2 logits lower than those with a +3.0 fitness
9 level. On the basis of the overall RMPFS measures, it would be possible to construct a norm based
10 scoring system for specific age groups if the interpretation of students' performance under a
11 norm-reference framework is needed. Although this study explored the theoretical possibility of
12 estimating an overall person measure of fitness based on very limited testing scores, it is not
13 recommended to determine person's fitness from only one assessment, because in practice the
14 information from just one indicator is too limited, i.e., the measurement error is too large, for practical
15 purposes. This overall measure combines core fitness components - cardiorespiratory fitness, muscular
16 endurance, and muscular strength - but is not the simple average of the performances on different
17 components. This simplified approach and reporting system provide a more efficient method and reduce
18 teachers' workload so that they could put more time and resources into the teaching and learning that
19 promotes children's physical fitness and health.

20 This research has its own limitations and future research which emphasises the following will
21 extend the contributions of this research to physical fitness assessment. The first, since this study took the
22 data fit the model approach, the indicators for two components of fitness (BMI for body composition and

1 Sit-and-Reach for flexibility) were excluded from the scale due to failure to fulfil the Rasch requirements
2 or practical considerations. Clearly body composition and flexibility are important and related both to
3 fitness and to health. Our conclusion is not that they are not important but our research reveals that they
4 do not behave, in the measurement sense, in the same way as do the other fitness indicators. At this stage,
5 our assertion is they should be assessed, recorded and utilised as indicators in their own right - but not
6 included in the construct of general fitness measure for these children. Future research could explore this
7 point further through two angles. One is to adhere to the “data fit the model” approach and to make efforts
8 to identify more appropriate indicators for the components of body composition and flexibility which can
9 be calibrated successfully into the Rasch measurement scale; those attempts could be made in smaller,
10 more closely controlled fitness testing contexts. The other is to explore multi-dimensional and continuous
11 Rasch models so as to identify a model with better fit the data. However, the model fits data approach is
12 likely to lose the strong measurement benefits which could be derived by adherence to Rasch model
13 principles.

14 Secondly, the RMPFS relies on the data exclusively from the partner school. That brings a
15 limitation to immediate generalization of the RMPFS developed in this research to application in other
16 samples. Thus this study’s core value remains in trying a new approach to physical fitness measurement
17 and building a good model practice - rather than providing a ready-for-use instrument for general physical
18 fitness assessment. Future research could utilize the techniques used in this research and extend to a larger
19 sample which might be representative for the whole Hong Kong primary or other school-aged student
20 population so that a Rasch measurement physical fitness scale could be developed for use with all Hong
21 Kong primary or other school-aged students. On the other hand, future research could use the same
22 technique to develop school-based databases for other similar samples to derive the same benefits. In

1 addition to replicating the practical benefits, there are also theoretical benefits that could be derived from
2 applying the same technique to other samples. The invariance of indicator measures (sample-distribution
3 free) required by Rasch model's means that indicator measures should be independent of any particular
4 sample used for indicator calibration. However, this research itself did not provide direct evidence of this
5 feature since it did not apply the RMPFS to other samples. Future investigations using already existing
6 data from other resources could provide further evidence of the validity to the RMPFS.

7 Finally, this research does not deny that psychometric approaches to data analysis, other than the
8 Rasch model, might be appropriate for producing more comprehensive descriptions of the variability in
9 this large longitudinal data set of children's physical fitness indicators. At the conclusion of this research,
10 it remains an open question as to whether other quantitative approaches might produce results that have
11 better fit of the model to the data. The completion of such an investigation could provide an interesting
12 complement to the results of the data fit the model Rasch measurement approach explicitly adopted at the
13 outset of this research.

14

1

Acknowledgment

2

This study is funded by a James Cook University Doctoral Publication Award. An earlier version of the

3

paper was presented to the Pacific Rim Objective Measurement Symposium (PROMS), at The Hong

4

Kong Institute of Education, Hong Kong, 28-30 July 2009. We thank the reviewers for their constructive

5

critique.

6

References

- 1
2 Andersen, E. B. (1995). What Georg Rasch would have thought about this book. In G. H. Fischer & I. W.
3 Molenaar (Eds.), *Rasch models: Foundations, recent developments, and applications* (pp. 383-390).
4 NY: Springer.
- 5 Andrich, D. (2004). Controversy and the Rasch model: A characteristic of incompatible paradigms?
6 *Medical Care, 42*, 1-16.
- 7 Bond, T. G., & Fox, C. M. (2007). *Applying the Rasch model: Fundamental measurement in the human*
8 *sciences* (2nd ed.). Mahwah, N.J.: Erlbaum.
- 9 Bowles, R. P., & Ram, N. (2006). Using Rasch measurement to investigate volleyball skills and inform
10 coaching. *Journal of Applied Measurement, 7*(1), 39-54.
- 11 B üsch, D., & Strauss, B. (2005). Qualitative differences in performing coordination tracks. *Measurement*
12 *in Physical Education and Exercise Science, 9*(3), 161-180.
- 13 B üsch, D., Strauss, B., Seidel, I., Pabst, J., Tietjens, M., Müller, L., Kretschmer, J., & Wirzsing, D. (2009).
14 Die Konstruktvalidit ä des Allgemeinen Sportmotorischen Tests für Kinder. *Sportwissenschaft, 39*,
15 95-103.
- 16 Fleishman, F. A. (1964). *The structure and measurement of physical fitness*. NJ: Prentice-Hall.
- 17 Hands, B., & Larkin, D. (2001). Using the Rasch measurement model to investigate the construct of
18 motor ability in young children. *Journal of Applied Measurement, 2*, 101 – 120.
- 19 Heesch, K. C., Masse, L. C., & Dunn, A. L. (2006). Using Rasch modeling to re-evaluate three scales
20 related to physical activity: Enjoyment, perceived benefits and perceived barriers. *Health Education*
21 *Research, 21*. Retrieved October 3, 2006, from <http://her.oxfordjournals.org/cgi/reprint/cyl054v1>
- 22 Heyward, V. H. (2002). *Advanced fitness assessment and exercise prescription* (4th ed.). Champaign, IL:
23 Human Kinetics.

-
- 1 Hong Kong Education and Manpower Bureau. (2005). *Hong Kong school physical fitness award schemes:*
2 *Teachers' handbook*. Retrieved August 10, 2006, from
3 <http://cd1.emb.hkedcity.net/cd/pe/tc/rr/pfas/handbook>
- 4 Hsueh, I. P., Wang, W. C., Sheu, C. F., & Hsieh, C. L. (2004). Rasch analyses of combining two indices
5 to assess comprehensive ADL function in stroke patients. *Stroke*, 35(3), 721-726.
- 6 Kang, J., & Kang, M. (2006). Rasch calibration of perceived weights of different sports games.
7 *Measurement in Physical Education and Exercise Science*, 10(1), 51-66.
- 8 Linacre, J. M. (2000). New approaches to determining reliability and validity. *Research Quarterly for*
9 *Exercise and Sport*, 71(2), 129-136.
- 10 Linacre, J. M. (2002). Optimizing rating scale category effectiveness. *Journal of Applied Measurement*, 3,
11 85-106.
- 12 Linacre, J. M. (2006a). *A user's guide to WINSTEPS/MINISTEPS: Rasch-model computer programs*.
13 Chicago, IL: Winsteps.com.
- 14 Linacre, J. M. (2006b). Data variance explained by Rasch measures. *Rasch Measurement Transactions*,
15 20(1), 1045.
- 16 Liu, Y. (2008). Youth fitness testing: If the "horse" is not dead, what should we do? *Measurement in*
17 *Physical Education and Exercise Science*, 12, 123–125.
- 18 Mahar, M. T., & Rowe, D. A. (2008). Practical guidelines for valid and reliable youth fitness testing.
19 *Measurement in Physical Education and Exercise Science*, 12(3), 126-145.
- 20 Marsh, H. W. (1993). The multidimensional structure of physical fitness: Invariance over gender and age.
21 *Research Quarterly for Exercise and Sport*, 64(3), 256-273.
- 22 Marsh, H. W., & Redmayne, R. S. (1994). A multidimensional physical self-concept and its relations to
23 multiple components of physical fitness. *Journal of Sport & Exercise Psychology*, 16(1), 43-55.

-
- 1 Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests* (Reprint, with
2 Foreword and Afterword by B. D. Wright, Chicago: University of Chicago Press, 1980).
3 Copenhagen, Denmark: Danmarks Paedagogiske Institut.
- 4 Silverman, S., Keating, X. D., & Phillips, S. R. (2008). A lasting impression: A pedagogical perspective
5 on youth fitness testing. *Measurement in Physical Education and Exercise Science*, 12(3), 146-166.
- 6 Smith, A. B., Rush, R., Fallowfield, L. J., Velikova, G., & Sharpe, M. (2008). Rasch fit statistics and
7 sample size considerations for polytomous data. *BMC Medical Research Methodology*, 8. Retrieved
8 July 7, 2009, from <http://www.biomedcentral.com/1471-2288/8/33>
- 9 Smith, E. V., Jr. (2001). Evidence for the reliability of measures and the validity of measure interpretation:
10 A Rasch measurement perspective. *Journal of Applied Measurement*, 2, 281-311.
- 11 Stratton, G., & Williams, C. A. (2007). Children and fitness testing. In E. M. Winter, A. M. Jones, R. C. R.
12 Davison, P. D. Bromley, & T. H. Mercer (Eds.), *Sport and exercise physiology testing guidelines* (pp.
13 211-223). NY: Routledge.
- 14 Strauss, B., B üsch, D. & Tenenbaum, G. (2007). New developments in measurement and testing. In G.
15 Tenenbaum & R. Eklund (Eds.), *Handbook of Sport Psychology* (3rd ed.) (pp. 737-756). Boston,
16 MA: Wiley.
- 17 Tenenbaum, G., Strauss, B., & B üsch, D. (2007). Applications of generalized Rasch models in Sport,
18 exercise and the motor domains. In M. v. Davier & C. Carstensen (Eds.), *Multivariate and Mixture*
19 *Distribution Rasch models* (pp. 347-357). New York: Springer.
- 20 Verhelst, N. D., & Glas, C. A. (1995). The one parameter logistic model. In G. H. Fischer & I. W.
21 Molenaar (Eds.), *Rasch models: Foundations, recent developments, and applications* (pp. 215-237).
22 NY: Springer.
- 23 Welk, G. J. (2008). The role of physical activity assessments for school-based physical activity promotion.
24 *Measurement in Physical Education and Exercise Science*, 12, 184–206.

-
- 1 Williams, C. S., Harageones, E. G., Johnson, D. J., & Smith, C. D. (2000). *Personal fitness: Looking*
2 *good / feeling good* (4th ed.). Dubuque, IA: Kendall/Hunt Publishing Company.
- 3 Wright, B. D., & Masters, G. N. (1982). *Rating scale analysis*. Chicago, IL: MESA Press.
- 4 Wright, B. D., & Mok, M. M. C. (2000). Rasch models overview. *Journal of Applied Measurement*, 1(1),
5 83-106.
- 6 Zhu, W. (2001). An empirical investigation of Rasch equating of motor function tasks. *Adapted Physical*
7 *Activity Quarterly*, 18(1), 72-89.
- 8 Zhu, W., & Cole, E. L. (1996). Many-faceted Rasch calibration of a gross motor instrument. *Research*
9 *Quarterly for Exercise and Sport*, 67(1), 24-34.
- 10 Zhu, W., & Kurz, K. A. (1994). Rasch partial credit analyses of gross motor competence. *Perceptual &*
11 *Motor Skills*, 79(2), 947-961.
- 12 Zhu, W., & Safrit, M. J. (1993). The calibration of a sit-up task using the Rasch Poisson Counts model.
13 *Canadian Journal of Applied Physiology*, 18(2), 207-219.
- 14 Zhu, W., Timm, G., & Ainsworth, B. (2001). Rasch calibration and optimal categorization of an
15 instrument measuring women's exercise perseverance and barriers. *Research Quarterly for Exercise*
16 *and Sport*, 72(2), 104-116.
- 17

1 **Tables and Figures**

2
3 Table 1

4 *Details of the Sample*

Academic year Semester	2002-03		2003-04		2004-05		2005-06		2006-07		Total
	1 st	2 nd	1 st	2 nd	1 st	2 nd	1 st	2 nd	1 st	2 nd	
Male	510	0	556	551	572	574	592	590	606	598	5149
Female	458	0	472	468	492	489	488	487	468	468	4290
Total	968	0	1028	1019	1064	1063	1080	1077	1074	1066	9439

5

Age	6	7	8	9	10	11	12	13	Missing	Total
Male	837	900	877	845	813	779	94	4	0	5149
Female	666	701	727	742	717	672	61	0	4	4290
Total	1503	1601	1604	1587	1530	1451	155	4	4	9439

6
7
8 Table 2

9 *Scale Properties of the RMPFS*

	Measure (logits)	S.E.	Infit MNSQ	Outfit MNSQ	Point-Measure Correlation
6-minute Run	-0.61	0.03	0.93	0.96	0.78
9-minute Run	1.25	0.03	0.85	0.88	0.86
1-minute Sit-ups	-1.59	0.02	0.95	1.00	0.79
Dominant Handgrip	0.96	0.02	1.11	1.13	0.79
Person	Separation:	1.83	Reliability:	0.77	
Item	Separation:	43.16	Reliability:	1.00	

10
11
12 Table 3

13 *Developmental Procedure for RMPFS*

	Infit MNSQ	Outfit MNSQ	Point-Measure Correlation	Rasch Reliability (Person/Item)	Variance Explained by Measures
8-indicator 9-category (Scale 1)					
6-minute Run	1.03	1.03	0.58		
9-minute Run	1.09	1.09	0.65		
1-minute Sit-ups	0.93	0.91	0.63		
Right Handgrip	0.76	0.75	0.73		
Left Handgrip	0.74	0.72	0.73	0.52/1.00	62.1%
Sit-and-Reach	1.21*	1.27*	0.42		
Standard Push-ups	1.01	1.01	0.70		
Modified Push-ups	1.49*	1.48*	0.47		
7-indicator 9-category (Scale 2)					
6-minute Run	1.10	1.10	0.61	0.66/1.00	60.6%

9-minute Run	1.15	1.15	0.68		
1-minute Sit-ups	1.07	1.05	0.64		
Right Handgrip	0.78	0.78	0.76		
Left Handgrip	0.75	0.75	0.76		
Standard Push-ups	1.01	1.02	0.74		
Modified Push-ups	1.57*	1.58*	0.51		
6-indicator 9-category (Scale 3)					
6-minute Run	0.93	0.92	0.70		
9-minute Run	0.95	0.95	0.75		
1-minute Sit-ups	0.90	0.90	0.69		
Dominant Handgrip	1.11	1.10	0.65	0.60/1.00	62.6%
Standard Push-ups	0.88	0.88	0.78		
Modified Push-ups	1.26*	1.27*	0.6		
4-indicator 9-category (Scale 4)					
6-minute Run	0.92	0.91	0.73		
9-minute Run	0.90	0.90	0.79		
1-minute Sit-ups	0.97	0.98	0.70	0.63/1.00	66.9%
Dominant Handgrip	1.09	1.08	0.70		
4-indicator 7-category (Scale 5)					
6-minute Run	0.92	0.95	0.72		
9-minute Run	0.90	0.91	0.79		
1-minute Sit-ups	0.95	0.99	0.73	0.62/1.00	68.7%
Dominant Handgrip	1.10	1.10	0.71		
4-indicator 7-category without Underfitting Persons (Scale 6 / RMPFS)					
6-minute Run	0.93	0.96	0.78		
9-minute Run	0.85	0.88	0.86		
1-minute Sit-ups	0.95	1.00	0.79	0.77/1.00	81.5%
Dominant Handgrip	1.11	1.13	0.79		

1 *Note: * Misfitting item*

2

3

4 Table 4

5 *Scale Properties of the RMPFS for Subsamples*

Sample	Measure (logits)	S.E.	Infit MNSQ	Outfit MNSQ	Point-Measure Correlation	Person Reliability /Separation	Item Reliability /Separation
Overall	R9	1.25	0.03	0.85	0.88	0.86	
	DH	0.96	0.02	1.11	1.13	0.79	
	R6	-0.61	0.03	0.93	0.96	0.78	0.77/1.83
	SU	-1.59	0.02	0.95	1.00	0.79	1.00/43.16
2003/1	R9	2.63	0.08	0.79	0.81	0.85	
	R6	-0.06	0.08	0.90	0.99	0.75	
	DH*	-1.13	0.06	1.12	1.12	0.78	0.78/1.87
	SU	-1.44	0.06	0.93	0.97	0.77	1.00/22.27

2004/1	R9	2.33	0.09	0.80	0.84	0.86	0.80/2.03	1.00/20.54
	DH	-0.18	0.06	0.99	1.00	0.82		
	R6	-0.31	0.08	1.01	1.04	0.72		
	SU	-1.84	0.06	0.98	1.05	0.76		
2005/1	R9	2.15	0.09	0.87	0.89	0.82	0.76/1.78	1.00/21.41
	DH	0.45	0.06	1.00	1.01	0.81		
	R6	-0.45	0.08	0.96	1.00	0.74		
	SU	-2.14	0.06	1.01	1.00	0.76		
2005/2	R9	3.37	0.08	0.74	0.75	0.87	0.80/2.02	1.00/27.82
	DH	-0.45	0.07	1.16	1.19	0.77		
	R6	-0.88	0.07	0.85	0.85	0.79		
	SU	-2.03	0.06	0.99	1.01	0.78		
2006/1	R9	1.56	0.09	0.82	0.83	0.86	0.76/1.77	1.00/18.89
	DH	0.93	0.06	1.11	1.15	0.74		
	R6	-0.46	0.08	0.99	1.04	0.69		
	SU	-2.03	0.06	0.92	1.01	0.74		
2006/2	R9	2.04	0.09	0.88	0.89	0.83	0.77/1.85	1.00/16.82
	DH	0.11	0.07	1.10	1.12	0.79		
	R6	-0.74	0.08	0.88	0.92	0.78		
	SU	-1.42	0.07	0.97	1.00	0.78		

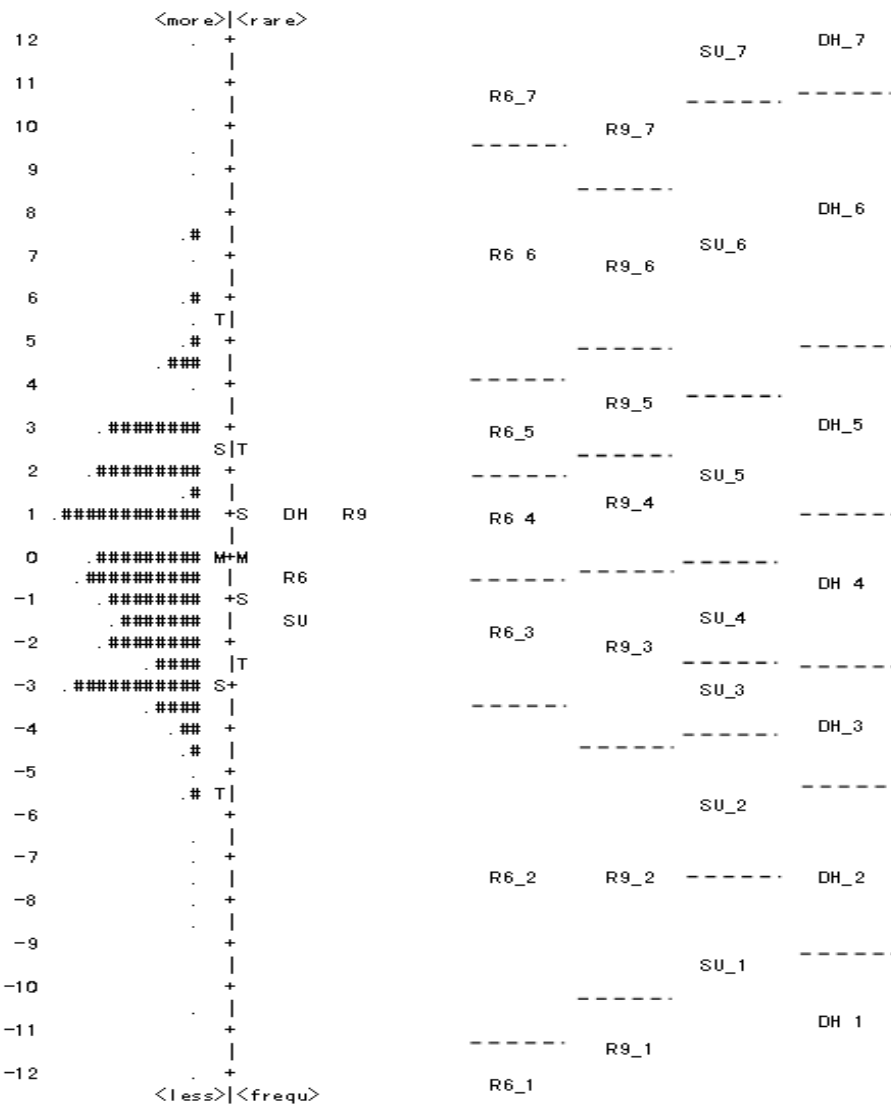
Note: *Indicator out of order.

Table 5

Category Equivalence for R6 and R9

Raw Scores (M)	Metres / min	R6	R9	Metres / min	Raw Scores (M)
Min- Max	Min- Max	Category		Min-Max	Min-Max
360-471	60-79	1		62-81	560-726
472-676	79-113	2		81-114	727-1027
677-809	113-135	3		114-136	1028-1221
810-969	135-162	4		136-161	1222-1452
970-1160	162-194	5		161-192	1453-1727
1161-1389	194-232	6		192-228	1728-2054
1390-1520	232-253	7		228-249	2055-2240

Note: The two columns at the left side present the metres covered and the speed for each category of R6; the two columns at the right side present the metres covered and the speed for each category of R9.

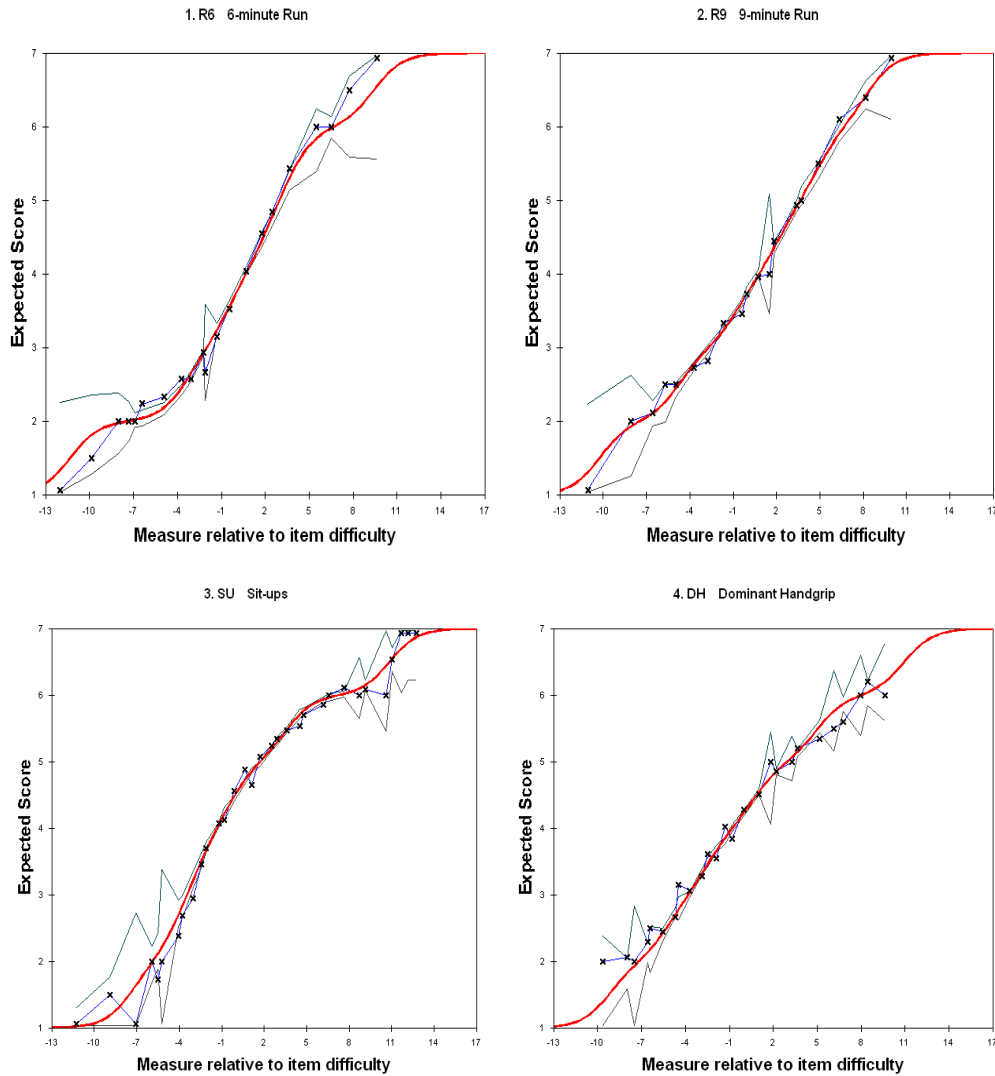


1

2 *Figure 1.* Wright map of the RMPFS.

3 Each '#' represents 64 children. R6: 6-minute Run; R9: 9-minute Run, SU: 1-minute Sit-ups; DH: Dominant

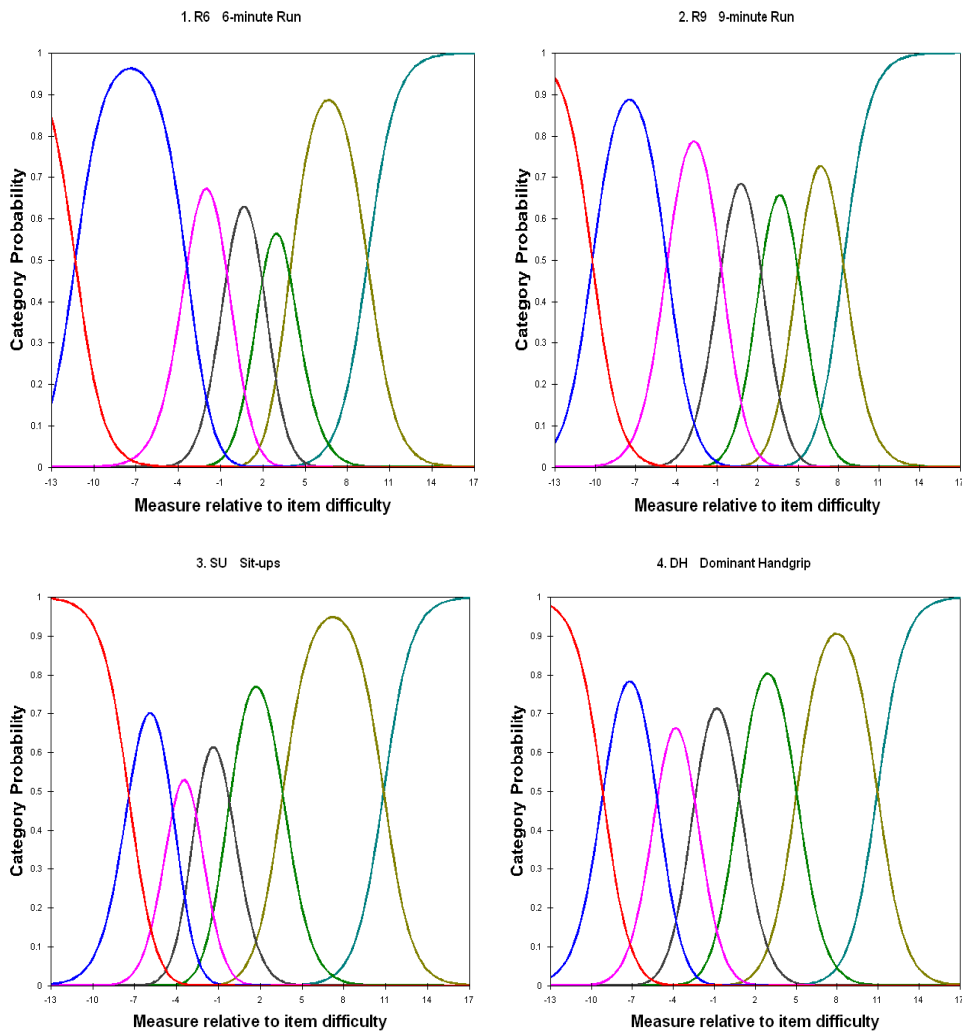
4 Handgrip.



1

2 *Figure 2.* Empirical (blue) and expected (red) ICCs for RMPFS indicators.

3 R6: 6-minute Run; R9: 9-minute Run, SU: 1-minute Sit-ups; DH: Dominant Handgrip.



1
2
3
4

Figure 3. Category probability curves for RMPFS indicators.

R6: 6-minute Run; R9: 9-minute Run, SU: 1-minute Sit-ups; DH: Dominant Handgrip.

1

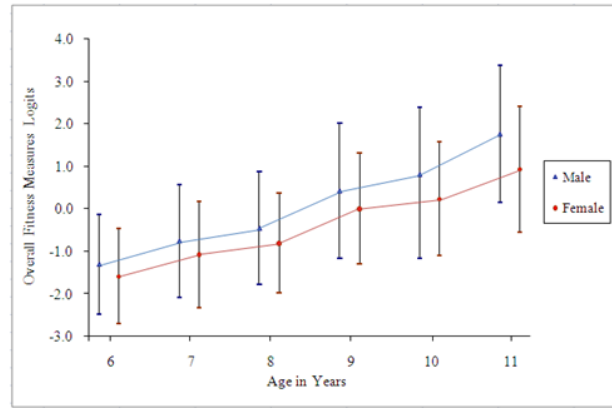


Figure 4a. Fitness development by age and sex (M±1S.D.)

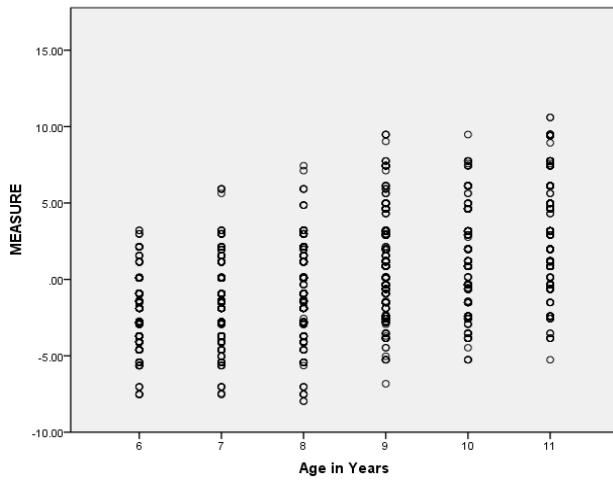


Figure 4b. Distribution of fitness levels for boys

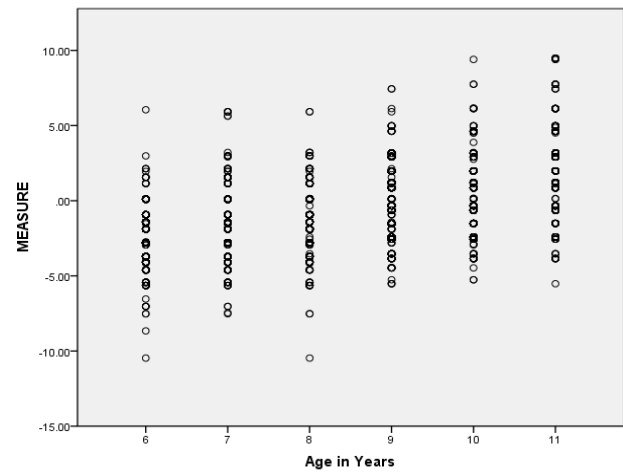


Figure 4c. Distribution of fitness levels for girls

2

3