

Inter-Rater Agreement and Test-Retest Reliability of the Performance and Fitness (PERF-FIT) Test Battery for Children: A Test for Motor Skill Related Fitness.

Bouwien Smits-Engelsman (✉ bouwiensmits@hotmail.com)

University of Cape Town Faculty of Health Sciences <https://orcid.org/0000-0003-0632-3276>

Eline Smit

Avans University of Applied Science: Avans Hogeschool Breda

Rosemary Xorlanyo Doe-Asinyo

UCT FHS: University of Cape Town Faculty of Health Sciences

Stella Elikplim Lawerteh

UCT FHS: University of Cape Town Faculty of Health Sciences

Wendy Aertssen

Avans University of Applied Science: Avans Hogeschool Breda

Gillian Ferguson

UCT FHS: University of Cape Town Faculty of Health Sciences

Dorothee L Jelsma

RUG: Rijksuniversiteit Groningen

Research article

Keywords: Agreement, Reliability, Physical fitness, Skill-related physical fitness, children, low resourced settings, Psychometric properties, Motor development

Posted Date: September 17th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-76118/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License. [Read Full License](#)

Version of Record: A version of this preprint was published on March 11th, 2021. See the published version at <https://doi.org/10.1186/s12887-021-02589-0>.

Abstract

Background: The Performance and Fitness (PERF-FIT) test battery for children is a recently developed, valid assessment tool for measuring motor skill-related physical fitness in 5 to 12-year-old children living in low-income settings. The aim of this study was to determine: (1) inter-rater agreement and (2) test-retest reliability of the PERF-FIT in children from 3 different countries (Ghana, South Africa and the Netherlands).

Method: For inter-rater reliability 29 children, (16 boys and 13 girls, 6-10 years) were scored by 2 raters simultaneously. For test-retest reliability 72 children, (33 boys and 39 girls, 5-12 years) performed the test twice, minimally one week and maximally two weeks apart. Relative and absolute reliability indices were calculated. ANOVA was used to examine differences between the three assessor teams in the three countries.

Results: The PERF-FIT demonstrated excellent inter-rater reliability (ICC, 0.99) and good test-retest reliability (ICC, ≥ 0.80) for 11 of the 12 tasks. A poor ICC was found for the Jumping item only, due to low spread in values. Overall, measurement error, Limits of Agreement and Coefficient of Variation had acceptable levels to support clinical use. No systematic differences were found between first and second measurement between the three countries but for one item (Overhead throw).

Conclusions: The PERF-FIT can reliably measure motor skill related fitness in 5 to 12-year-old children in different settings and help clinicians monitor levels of power and agility, and fundamental motor skills (throwing, bouncing, catching, jumping, hopping and balance).

Introduction

Despite the global interest in promoting physical activity and fitness among school-aged children, there is a paucity of studies concerning this topic from developing countries. The few studies available provide data that children living in socioeconomically disadvantaged circumstances are disproportionately experiencing limited opportunities to develop adequate levels of physical fitness (Valentini, Clark, Whittall, 2015). The levels of physical fitness and motor skills in children are important factors to be able to participate in daily activities such as sports. Child factors can partly determine the level of physical fitness or motor skills, but the environment can be as influential (Ortega, Ruiz, Castillo, Sjöström, 2008). It is reported that South African children living in low income areas have overall lower levels of aerobic fitness compared to children in other settings (Aertssen, Bonney, Ferguson, & Smits-Engelsman, 2018). If children experience difficulties performing motor skills in daily life, early identification is important for timely intervention purposes. However, no valid and reliable assessment tools existed for measuring motor skill-related physical fitness (Table 1) in young children across low resource areas (Caspersen, Powell, Christenson, 1985). Given that children in low resource areas are largely underrepresented in motor development research, it was prudent to develop a new test called Performance and Fitness (PERF-FIT) test battery for this target population that evaluates motor skill-related physical fitness and integrates muscular fitness and motor coordination, as this is deemed more ecologically valid. The test was developed for health and teaching professionals (occupational and physical therapist, school nurses and physical educators). Consistent with the World Health Organization's (WHO, 2007) International Classification of Functioning, Health and Disability, the PERF-FIT is developed as a tool measuring the "activity" component of the WHO framework rather than the "body structure and function". The rationale for this focus is the desire to detect deviations in the development of motor skills and fitness levels that have a functional impact on the day-to-day activities of children. The benefit of such a tool is early identification of children with deficits in fundamental movement skills and muscular fitness. The detection of such variations will enable researchers and clinicians to explore possible etiological mechanisms and policy makers to develop preventive measures. If needed this could lead to clinical actions.

After having established good content and structural validity of the PERF-FIT in Brazil (Smits-Engelsman, Bonney, Neto, & Jelsma, 2020; Smits-Engelsman, Cavalcante Neto, Draghi, Rohr, & Jelsma, 2020), this study examined how consistent scores on the PERF-FIT are under different circumstances and different populations. The purpose of the first part of the study was to check if the item instructions for scoring, would lead to comparable results between two raters, when assessing the performance of the child at the same time (inter-rater agreement).

Next, we evaluated test-retest reliability. Usually, clinical measurement is rarely perfectly reliable as raters and subjects are known to respond with some inconsistency. Since reliability is generally population specific, a comparison of reliability between different populations is advised (Bruton, Conway, & Holgate, 2000). Due to expected use of the PERF-FIT in very different contexts, we collected data in three countries with different groups of raters. Subjects of three different populations were tested twice, in order to test the stability of the measure over time.

Method

Participants

In total, 101 children between 5–12 years of age were recruited in two elementary schools in low income areas in Cape Town, South Africa (SA), two elementary schools in low income areas in Accra, Ghana (GH) and in two elementary schools in middle income areas in Tilburg, the Netherlands (NL) (See Fig. 1). The sample was randomly chosen by the teachers from the class list.

First, we examined the inter-rater agreement of the PERF-FIT test battery. Twenty-nine 6-10-year-old South African children were included in this part of the study. Next, we examined test-retest reliability to evaluate possible variance in performance between two test moments in the children and if the variance was stable under the different testing circumstances. In total, 72 children between 5 and 12 years of age completed the test-retest part of the study; South African children (24), Ghanaian children (23) and Dutch children (25). The ethical review committees of the University of Cape Town, Ghana Health Service and University of Groningen gave their approval for the study (UCT HREC Ref 598/2019; HREC139/2019; GHS-ERC 084/04/19; PSY-1920-S-0107). Demographic characteristics are summarized in Fig. 1.

Insert Fig. 1 about here

Procedure

Permission to approach the head teachers were obtained from the school districts. Verbal and/or written explanations of study purpose, test procedures, benefits and risks were provided to parents. Children were included after parents or caretakers signed the consent forms and children gave assent to participate. Children included were a random sample of school children aged 5–12 years and with understanding of the local language. Children with health-related conditions were excluded based on the Physical Activity Readiness Questionnaire (PAR-Q) (Warburton, et al., 2011). In addition to PERF-FIT scores, data sought included age, height, weight and gender. No other information was available to the raters about the children.

Assessments were performed under the circumstances that the test was developed to be used in; on the school's premises outside (GH and SA), in the gym or hall (NL and SA) and a physiotherapy practice (NL). Participants completed standardized warm-up procedures prior to testing as prescribed in the manual. They were allowed practice trials for each test item before the scored trial as indicated in the manual. Children who did not have proper shoes, performed the test barefoot on both occasions. All of the Ghanaian children were tested barefoot; part of the South African children wore uniform shoes and part was barefoot. All the Dutch children wore sneakers.

The lead author trained at least one rater per country but was not present at any of the test sessions. The trained raters instructed the other raters in SA, GH and NL during a half-day training, where they practiced in small groups to obtain a solid routine. Raters were selected as being representative of the future users; pediatric physiotherapists, physiotherapists and occupational therapist, teaching assistants and a nurse. The raters conducted all the testing during school hours except in the Netherlands where part of the testing was done on a day-off.

Inter-rater Agreement Study

The consistency of two different clinicians rating the PERF-FIT was tested. When establishing inter-rater agreement with two observers, one tests if the instructions for scoring were unambiguous and if this led to similar results. Overall results are excellent (mean ICC 0.99), indicating that the two raters did get the same results for the same subjects. Since the children were selected randomly by the teachers, the results can be generalized for the child population within this age range (D'Olhaberriague et al., 1996).

Test-retest Reliability Study

Test-retest reliability concerns the reproducibility of the observed value when the measurement is repeated in a stable population. Studying reliability may seem straightforward, as it is just a matter of repeating the measurement on a reasonable number of individuals. However, interpreting the findings is less simple and a combination of approaches is more likely to give a true picture of reliability (Bruton, et al., 2000). The type of data (continuous) of the PERF-FIT requires standard error of measurement (SEM) (De Vet et al., 2006; Stratford and Goldsmith, 1997) and proportions of agreement within specified limits to provide useful information concerning reliability (De Vet et al., 2006). Given the ICC's found in this study, one can assume that the PERF-FIT is a reliable tool. ICC's for 4 items are 80 or higher and 7 items have an ICC of 90 or higher. The relative nature of the ICC is reflected in the fact that the magnitude of an ICC depends on the between-subjects variability. That is, if subjects differ little from each other (homogeneous sample), ICC values can be low even if trial-to-trial variability is small as shown in the *Jumping* item. This item, which is easy in this population, showed low ICC but good agreement (85%). It would also be of interest to test the reliability of this item in young children and with DCD. Importantly, if we were to include participants with neurodevelopmental delays the between-subjects variability will change as well as the ICC (Strainer, Norman, & Cairney, 2014).

ICC is not sensitive to disagreement due to systematic bias as was shown by the comparison between test 1 and test 2, half the items showed a very small but significant improvement but have high ICC's. The need to perform the test twice will cause performance variability, due to changes in motivation and familiarization with the tasks. Detailed analysis of the *Side jump* data, with good ICC (0.90), showed that five children "improved" ten jumps or more (max 13). However, this was not due to instruction or circumstances since the five children came from three different countries. Still these differences cannot be attributed to improved anaerobic fitness, or improved motivation since these children showed no improvement on the other items. Hence this finding points more towards a short-term learning effect or getting the clue of the agility required in this task for some children. We therefore added the recommendation in the manual to offer one extra practice opportunity if a child is still struggling with the weight shifts of the Side jump.

Throwing and catching series also showed small improvements. Some of the African children were less used to this task, which may have increased the learning effect. Consequently, we will emphasize to consistently use the two practice trials *per level of difficulty*, to reduce the learning effect.

Test location. The subject population of interest for the PERF-FIT is the group of children in elementary school age living in low socio-economic circumstances. Children with different lifestyles (level of daily physical activity, participation in structured physical education and sports) and testing in different contexts may respond differently to re-testing of some tasks. Therefore, we gathered data in three countries with many raters ($n = 16$), to analyze the reliability across these different populations and environments making the results clinically more widely applicable (Bruton, Conway, & Holgate, 2000). Although the testing was done in a standardized way, raters, sites and children were very different. Still, no country-related bias was found except for the *Overhead throw*, where the difference in scores between the two test occasions was larger in the Ghanaian children. Scoring this item requires the tester to focus on the landing spot, preferably on sand, dirt floor or grass so the sandbag leaves a landing imprint. The practice trial given in this task is done with submaximal force, which may have decreased the familiarization in the first testing.

Despite the noise and distraction, inherent to testing at the school premises in open space, the test results were considerably stable, which implies that the children were able to attend to the instructions under these circumstances. These findings point to the fact this test is enjoyable and engaging for the children. It is to be expected that if children are tested in a more clinical one-on-one situation, the variability between test and retest will be even less.

In this study we choose for a wide variety of outcomes because they all have advantages and disadvantages. Both the SEM and LoA were calculated because they differ in the type of measurement error that they describe and in the coverage probability of the reference interval (0.68 versus 0.95%). If the variability in test-retest outcomes depends on the magnitude of the mean values, the use of a ratio statistic is useful to the researchers. The advantage of CV being unitless is that it can be used to compare different instruments, but this makes it harder to translate results into clinical practice.

Outcome Measure: Perf-fit

The PERF-FIT measures motor skill related physical fitness in children aged between 5–12 years old. The test has two subscales: a Performance part and a Fitness part. See Table 1. The PERF-FIT test battery is easy to administer, low-cost and developed for measuring performance-related physical fitness in school-aged children living in low-income settings and has excellent content validity and good structural validity (Smits-Engelsman et al., 2020a, 2020b). A full description of the PERF-FIT test battery is available through the first author (Smits-Engelsman 2018).

Table 1
Items of the PERF-FIT

PERF-FIT	
<i>Motor Skill Performance items</i>	
Bouncing and Catching	Children bounce tennis ball to the floor and catch. This series involves five bouncing and catching items of increasing skill difficulty. All children start at the easiest level. This series is discontinued if the child scores less than 6 out of 10 catches.
Throwing and Catching	Children throw tennis ball in the air to at least eye level height and catch. This series involves five throwing and catching items of increasing skill difficulty. All children start at the easiest level of this series. The series is discontinued if the child scores less than 6 out of 10 catches.
Jump	Children are asked to jump inside an agility ladder. This series involves four jumping items of increasing difficulty. Two test trials are allowed if maximum score is not obtained.
Hop	Children are asked to hop inside an agility ladder. This series involves four hopping items of increasing difficulty for each leg. Two test trials are allowed if maximum score is not obtained.
Balance	Children are asked to perform two (2) static balance tasks for each leg and three (3) dynamic balance tasks. Tasks involve knee hugging, grasping the foot and picking cans from the floor at close and far distance.
<i>Agility and Power items</i>	
Running	Children are asked to run (one foot per square) in 3.5 m agility ladder and run around a bottle placed at a distance of 50 cm from the starting line and return the same way as fast as possible. Two test trials are given for each child. The time taken (in seconds) to complete this task and number of mistakes made are recorded.
Stepping	Children are made to step with two feet in each square of a 3.5 m agility ladder and run around a bottle placed at a distance of 50 cm from the starting line and return the same place as fast as possible. Two test trials are given for each child. The time taken (in seconds) to complete this task and number of mistakes made are recorded.
Side Jump	Children are required to jump sideways on their feet. One foot per square, in the same three squares of the agility ladder. The total number of correct landings in 15 s is recorded for each of the two test trials.
Long Jump	Children are asked to jump forward as far as possible and land on their feet in a controlled manner (i.e. balanced landing). The distance between the starting line and the heel of the foot that landed closest to the starting line is measured in centimeters. Two test trials are given.
Overhead Throw	Children kneel just behind a starting line and throw a sandbag (2 kg) forward as far as possible. The bag is held over the head and thrown from a starting position behind the head. The distance between the starting line and the part of the sandbag closest to the starting line is measured in centimeters. Two test trials are performed.

Agility and power subscale

This subscale contains five items: *Running*, *Stepping*, *Side Jump*, *Long Jump*, and *Overhead Throw*. For the *Agility and power subscale* children perform two trials for each item and get 15 seconds rest in between.

Motor Skill Performance subscale

This subscale contains five Skill Item Series (SIS) of increasing difficulty; *Bouncing and Catching*, *Throwing and Catching*, *Jumping*, *Hopping* (left and right), and *Balance*. All children start at the easiest level and a series is terminated when they do not reach the criterion number of points for the item after two trials. If a child obtains the maximum number of points after the first trial no second trial is given and the child proceeds to the next level of difficulty.

After the first round of collecting validity data in Brazil (Smits-Engelsman et al., 2020a), it was found that most children obtained a maximum score on the static balance series and it was decided to increase the total number of seconds of the static balance series from 40 to 60 seconds for future studies. At this moment the data collection for SA had already started with the 40 seconds protocol. Therefore, the South African data on one item, *Static balance*, was discarded in the current paper. This was the only adaptation in the protocol, which then was used for data collection in GH and NL.

Data analysis

Descriptive data were calculated in terms of mean value and standard deviation (Mean \pm SD). Relative reliability, which is the degree to which individuals maintain their position in a sample over repeated scoring or testing, was determined by calculating the two-way random intra-class correlation coefficient (ICC 2,1a) for absolute agreement of single measures. The 95% confidence interval (CI) was calculated for each ICC (Shrout, & Fleiss, 1979). Reliability was considered poor for ICC values < 0.40 , fair for values between $0.40-0.59$, good for values between

0.60–0.74, and excellent for values between 0.75–1.00 (Cicchetti, 1994; Cicchetti et al., 2006). ICC values above 0.75 were considered acceptable for test-retest reliability. (Portney, & Watkins, 2009).

A paired t-test was used to compare the means of test (T1) and retest (T2) to evaluate whether there was any statistically significant bias between the test results.

Next, indicators of absolute reliability were calculated to determine the degree to which repeated measurements vary for individuals, expressed in the actual units of measurement, or as a proportion of the measured values. The Standard Error of Measurement (SEM), Bland and Altman's 95% Limits of Agreement (LoA) (Bland and Altman 1986) and coefficient of variation (CV) are all measures of absolute reliability that were used in this study.

The calculation of SEM and LoA do not depend on sample size, but the precision of their estimate for the population parameter does. Bland and Altman recommended sample sizes of at least 50 individuals in a study to consider the sample LoA to be a precise estimate of the population LoA (Bland, & Altman, 1986). Since we were also interested in a group comparison we oversampled, and we aimed at 25 subjects per country.

The SEM, as measure of precision of the assessment, was determined using the ICC through the formula $SEM_{\text{agreement}} = SD \cdot \sqrt{(1 - ICC_{\text{agreement}})}$ in which SD is the sample SD of the grand mean and ICC is the calculated intraclass correlation coefficient (Weir, 2005).

Minimal Detectable Change (MDC) was calculated as $MDC_{95} = 1.96 \cdot \sqrt{2} \cdot SEM_{\text{agreement}}$ (de Vet, Terwee, Mokkink, & Knol, 2011; Haley & Fragala-Pinkham, 2006). The MDC_{95} is the minimal amount of change observed before the change can be considered to exceed the variation and measurement error at the 95% confidence level.

Absolute reliability statistics were also calculated using the standard deviation of test-retest differences ($SD_{\text{differences}}$) and its derivatives. $SD_{\text{differences}}$ is the SD of the differences between values on T1 and T2.

The 95% LoA were calculated as the mean difference $\pm (1.96 \cdot SD_{\text{differences}})$ (Blant Altman 1986; de Vet, Terwee, Mokkink, & Knol, 2011).

The Coefficient of Variation (CV) or relative standard deviation is the individual SD expressed as a percent of the mean of T1 and T2 using the formula $(SD/\text{Mean}) \cdot 100$. The higher the SD, the greater the percentage of the mean and the higher the %CV. A %CV of < 10% is considered excellent, 10–20% medium, implying good precision, 20–30% high, meaning low precision and > 30% is considered very high, indicating very low precision (Atkinson, & Nevill, 1998; Lee, et al., 2013).

To test for possible dissimilarities in the degree of the error between T1 and T2 in the participating countries an ANOVA was run on the difference score (T1-T2) for all items with country (3) as between group factor and post hoc Bonferroni tests.

Statistical data analyses were carried out using SPSS version 25.0. A value of $p < .05$ was considered statistically significant in all analyses.

Results

Inter-rater agreement

Very high ICC's were found ≥ 0.98 for all items. The results of the inter-rater agreement ($n = 29$) of the two raters are shown in Table 2.

Table 2
Inter-rater agreement of the PERF-FIT with Interclass Correlation Coefficient (ICC) and 95% confidence interval (CI) per item.

		Inter-rater agreement		
		ICC	95% CI Low	95% CI High
1	Running (s)	0.995	0.990	0.998
2	Stepping (s)	0.980	0.945	0.991
3	Side jump (#)	0.997	0.995	0.999
4	Long jump (cm)	1.00	1.00	1.00
5	Overhead throw (cm)	0.997	0.993	0.999
6	Bounce and Catch (#)	0.999	0.998	0.999
7	Throw and Catch (#)	1.00	1.00	1.00
8	Jump (#)	0.993	0.992	0.998
9a	Hop Left (#)	0.993	0.986	0.997
9b	Hop Right (#)	0.997	0.995	0.999
10a	Static balance (s)	0.986	0.971	0.994
10b	Dynamic balance (#)	0.994	0.988	0.997
S: items measured in seconds; cm: measured in centimeters; #: measured in number of times				

Test-retest Reliability

Test-retest reliability results (n = 72) of the sixteen raters in the three countries are depicted in Table 3.

Table 3

Test-retest reliability outcomes of the PERF-FIT item scores. Means per test occasion (Time 1 and Time 2) and Grand mean, Intraclass Correlation Coefficient (ICC) with 95% Confidence Interval, Standard Error of Measurement (SEM), Minimal Detectable Change (MDC₉₅), Mean difference (Mean Dif) between test occasion (SD), Limits of Agreement (LoA) with upper and lower limit, percentage Coefficient of Variation (CV), p-values for the t-test comparison between T1 and T2.

PERF-FIT items	Time 1	Time 2	Grand Mean	ICC Test-retest	ICC 95% CI	SEM	MDC ₉₅	Mean Dif 1-2	SD Dif	Lower limit	LoA	Upper limit	CV	p-value
Item 1 Running (s)	7.14	7.22	7.18	0.82	0.71–0.89	0.50	1.39	-0.08	1.01	-2.06	1.98	1.90	7.3	0.52
Item 2 Stepping (s)	14.19	13.29	13.74	0.80	0.64–0.89	1.09	3.01	0.90	2.01	-3.04	3.94	4.85	8.5	< 0.01
Item 3 Side jump (#)	25.3	28.6	26.9	0.90	0.57–0.96	2.5	6.81	-3.3	3.81	-10.77	7.46	4.16	9.5	< 0.001
Item 4 Long jump (cm)	120.1	120.6	120.4	0.90	0.84–0.94	7.3	20.36	-0.4	14.69	-29.22	28.79	28.36	6.3	0.80
Item 5 Overhead throw (cm)	212.2	212.9	212.5	0.95	0.92–0.97	13.0	36.17	-0.7	27.06	-53.75	53.04	52.33	6.9	0.83
Item 6 Throw and catch (max 50#)	36.4	38.2	37.3	0.96	0.92–0.97	2.7	7.37	-1.8	5.41	-12.42	10.60	8.78	13.1	0.01
Item 7 Bounce and catch (max 50#)	39.0	43.1	41.1	0.92	0.73–0.96	3.4	9.33	-4.1	5.79	-15.43	11.34	7.26	11.2	< 0.001
Item 8 Jump (max 20#)*	19.1	19.6	19.3	0.36	0.01–0.59	0.8	2.35	-0.5	1.67	-3.80	3.27	2.74	3.5	0.01
Item 9b Hop Left (max 20#)	13.9	14.6	14.3	0.90	0.84–0.94	2.0	5.62	-0.7	3.99	-8.54	7.81	7.09	21.0	0.13
Item 9a Hop Right (max 20#)	15.0	16.0	15.5	0.92	0.86–0.95	1.6	4.37	-1.1	3.12	-7.16	6.11	5.05	13.6	0.01
Item 10a Static balance (max 60 s)§	54.9	55.8	55.3	0.88	0.79–0.93	3.5	9.70	-1.0	7.10	-14.89	13.92	12.94	5.3	0.34
Item 10b Dynamic balance (max 32#)	26.2	26.6	26.4	0.88	0.81–0.92	1.9	5.34	-0.4	3.88	-8.05	7.61	7.16	8.6	0.34

PERF-FIT items	Time 1	Time 2	Grand Mean	ICC Test-retest	ICC 95% CI	SEM	MDC ₉₅	Mean Dif 1-2	SD Dif	Lower limit	LoA	Upper limit	CV	p-value
*% Agreement for score +/- 1 point = 84.7%. \$ Data from Ghana and the Netherlands. Max: maximum score. S: item measured in seconds; cm: item measured in centimeters; #: measured in number														

Overall test-retest reliability was good to excellent on 11 of the 12 items; all ICC's were .80 or higher (Table 3). Only the item *Jumping* showed a low ICC due to lack of spread in the data. This was the easiest item and many children had a maximum score (63% and 78% in T1 and T2, respectively). Percentage agreement plus or minus 1 point was 84.7%.

Comparison between first and second test occasion showed that there was a statistically significant difference on half of the items. However, as shown in Table 3 (Column Mean Difference 1-2) this systematic bias was small, except for *Bouncing and catch* and *Side jump* ($p < 0.001$). The SD of the differences in scores between the two test occasions and LoA for each variable with its 95% confidence interval are also shown in Table 3. The mean %CV is 9.6%, which indicates excellent stability and the highest %CV (21% for *Hopping* on left foot) was still considered acceptable.

Comparison Per Country

The repeated measure ANOVA showed that the differences between T1 and T2 were not significantly different between countries for 11 of the 12 scores (Table 3). Only the *Overhead throw* the difference was larger in the Ghanaian children. Post hoc test showed that they were different from the Dutch children, who were slightly worse on the second test while the Ghanaian children in general performed better the second time on this item (see Fig. 2).

Insert Fig. 2 About here

Discussion

A new tool, the PERF-FIT was developed because none of the currently available norm referenced motor performance tests for children of elementary school age combined fundamental skills and muscular skeletal fitness. This study aimed to evaluate whether the PERF-FIT is a reliable tool and whether the measurement error is acceptable for practical use. Because widely accepted criteria or guidelines for reliability and agreement reporting in the health care and medical fields are lacking (Kottner et al., 2011), we chose for a wide variety of outcomes to evaluate the reliability of the PERF-FIT. Inter-rater agreement depends primarily on good training of the raters, and on good standardization and description of the tasks (Smits-Engelsman, 2018). Data in this study indicate that this was the case for the PERF-FIT. Test-retest reliability is highly dependent on the situation or on the state and stability of the participants, and is therefore characterized by larger variability, which was confirmed by our results although the agreement between the first and second test occasion was good. A small learning or familiarization effect was seen in 6 of the 12 items. No systematic differences between test-retest differences were found between the testing sites in the three countries in randomly selected children between 5-12 years old, except for 1 item. An average CV of 10% - obtained in the current study- means that, assuming the data are normally distributed, 68% of the differences between tests lie within 10% of the mean of the data (Atkinson & Neville, 1998).

Inter-rater agreement study

The consistency of two different clinicians rating the PERF-FIT was tested. When establishing inter-rater agreement with two observers, one tests if the instructions for scoring were unambiguous and if this led to similar results. Overall results are excellent (mean ICC 0.99), indicating that the two raters did get the same results for the same subjects. Since the children were selected randomly by the teachers, the results can be generalized for the child population within this age range (D'Olhaberriague et al., 1996).

Test-retest reliability study

Test-retest reliability concerns the reproducibility of the observed value when the measurement is repeated in a stable population. Studying reliability may seem straightforward, as it is just a matter of repeating the measurement on a reasonable number of individuals. However, interpreting the findings is less simple and a combination of approaches is more likely to give a true picture of reliability (Bruton, et al., 2000). The type of data (continuous) of the PERF-FIT requires standard error of measurement (SEM) (De Vet et al., 2006; Stratford and Goldsmith, 1997) and proportions of agreement within specified limits to provide useful information concerning reliability (De Vet et al., 2006). Given the

ICC's found in this study, one can assume that the PERF-FIT is a reliable tool. ICC's for 4 items are 80 or higher and 7 items have an ICC of 90 or higher. The relative nature of the ICC is reflected in the fact that the magnitude of an ICC depends on the between-subjects variability. That is, if subjects differ little from each other (homogeneous sample), ICC values can be low even if trial-to-trial variability is small as shown in the *Jumping* item. This item, which is easy in this population, showed low ICC but good agreement (85%). It would also be of interest to test the reliability of this item in young children and with DCD. Importantly, if we were to include participants with neurodevelopmental delays the between-subjects variability will change as well as the ICC (Strainer, Norman, & Cairney, 2014).

ICC is not sensitive to disagreement due to systematic bias as was shown by the comparison between test 1 and test 2, half the items showed a very small but significant improvement but have high ICC's. The need to perform the test twice will cause performance variability, due to changes in motivation and familiarization with the tasks. Detailed analysis of the *Side jump* data, with good ICC (0.90), showed that five children "improved" ten jumps or more (max 13). However, this was not due to instruction or circumstances since the five children came from three different countries. Still these differences cannot be attributed to improved anaerobic fitness, or improved motivation since these children showed no improvement on the other items. Hence this finding points more towards a short-term learning effect or getting the clue of the agility required in this task for some children. We therefore added the recommendation in the manual to offer one extra practice opportunity if a child is still struggling with the weight shifts of the Side jump.

Throwing and catching series also showed small improvements. Some of the African children were less used to this task which may have increased the learning effect. Consequently, we will emphasize to consistently use the two practice trials *per level of difficulty*, to reduce the learning effect.

Test location. The subject population of interest for the PERF-FIT is the group of children in elementary school age living in low socio-economic circumstances. Children with different lifestyles (level of daily physical activity, participation in structured physical education and sports) and testing in different contexts may respond differently to re-testing of some tasks. Therefore, we gathered data in three countries with many raters (n=16), to analyze the reliability across these different populations and environments making the results clinically more widely applicable (Bruton, Conway, & Holgate, 2000). Although the testing was done in a standardized way, raters, sites and children were very different. Still, no country-related bias was found except for the *Overhead throw*, where the difference in scores between the two test occasions was larger in the Ghanaian children. Scoring this item requires the tester to focus on the landing spot, preferably on sand, dirt floor or grass so the sandbag leaves a landing imprint. The practice trial given in this task is done with submaximal force, which may have decreased the familiarization in the first testing.

Despite the noise and distraction, inherent to testing at the school premises in open space, the test results were considerably stable, which implies that the children were able to attend to the instructions under these circumstances. These findings point to the fact this test is enjoyable and engaging for the children. It is to be expected that if children are tested in a more clinical one-on-one situation, the variability between test and retest will be even less.

In this study we choose for a wide variety of outcomes because they all have advantages and disadvantages. Both the SEM and LoA were calculated because they differ in the type of measurement error that they describe and in the coverage probability of the reference interval (0.68 versus 0.95%). If the variability in test-retest outcomes depends on the magnitude of the mean values, the use of a ratio statistic is useful to the researchers. The advantage of CV being unitless is that it can be used to compare different instruments, but this makes it harder to translate results into clinical practice.

Limitations And Future Research

Given the way the inter-rater reliability was examined, variability as a result of instruction was not tested. During field-based testing, all sources of variability cannot be controlled, therefore the design chosen for this study is close to the context this test was developed for. Results of agreement and reliability studies are intended to provide information about the amount of error inherent to a measurement tool in a specific population and context. High ICC's reflect adequate relative reliability for use of the PERF-FIT in the population that has been investigated. However, measures of reliability are generated by distribution-based methods and are dependent on the mean and variance in the group. The Minimal Detectable Change is very susceptible to increased variance given its formula. Reliability studies should be repeated in the population the instrument will be applied in, since variability may be different in groups on children with known poor motor performance, low levels of fitness, or learning disabilities. Also, the impact of BMI on the scores and the reliability should be investigated in different weight categories. Additionally, studies are needed to evaluate the responsiveness of the PERF-FIT or ability of the test to measure changes after intervention.

Conclusion

The present study examined inter-rater agreement and test-retest reliability of the PERF-FIT in a manner that replicates how the test is typically used in the actual everyday context. Inter-rater agreement and test-retest reliability were adequate to support clinical use. Hence, the PERF-FIT was relatively stable over time based on the small differences between the repeated measurements and based on the calculated SEM's. The Coefficient of Variation on average was 10%, indicating good stability. Hardly any systematic differences were found between the testing sites in the three countries, which supports the use of the PERF-FIT by trained raters from a variety of backgrounds in different contexts.

Abbreviations

PERF-FIT: Performance and Fitness test battery

ICC: Intraclass Correlation Coefficient

SA: South Africa

GH: Ghana

NL: Netherlands

PAR-Q: Physical Activity Readiness Questionnaire

SEM : Standard Error of Measurement

LoA: Limits of Agreement

CV: Coefficient of Variance

SD: Standard Deviation

SIS: Skills Item Series

MDC: Minimal Detectable Change

ANOVA: Analysis of Variance

DCD: Developmental Coordination Disorder

BMI: Body Mass Index

Declarations

Contributors All individuals listed as authors meet the appropriate authorship criteria and have approved the acknowledgement of their contributions. The primary author, BCM, was responsible setting up the project, development of research design, for the drafting of the paper and liaising with the coauthors on findings and conclusions. DJ contributed to the paper through interpretation of data, completing methodological assessments and revising manuscript content throughout its development. JF was responsible for the logistics of the whole project and rater supervision. WA supervised data collection in the Netherlands, RD and SL were responsible for the project in Ghana, JF, ES and DJ for the project in SA. All contributed to the paper through assisting with the interpretation of data and revising manuscript content through its development.

Funding No funding

Competing interests The authors declare that they have no competing interests.

Consent for publication Not applicable.

Patient consent Ethical approval was obtained from the University of Cape Town, Ghana Health Service and University of Groningen gave their approval for the study (UCT HREC Ref 598/2019; HREC139/2019; GHS-ERC 084/04/19; PSY-1920-S-0107). Written informed consent to participate was obtained from the parents/guardians of the minors included in this study and assent was signed by the children. Permission was also obtained from the head teachers of the schools.

Data sharing statement The datasets used and analyzed during the current study are available from the corresponding author on reasonable request. The PERF-FIT manual and instruction videos can be accessed free of charge for the intended users after registration via the first author for use in low resource communities.

Acknowledgement We acknowledge the support of parents, children and management of the participating schools.

References

1. Aertssen W, Bonney E, Ferguson G, Smits-Engelsman B. Subtyping children with developmental coordination disorder based on physical fitness outcomes. *Human Movement Science*. 2018;60:87–97. <https://doi.org/10.1016/j.humov.2018.05.012>.
2. Atkinson G, Nevill AM. Statistical methods for assessing measurement error (reliability) in variables relevant to sports medicine. *Sports Med*. 1998;26:217–38.
3. Bland JM, Altman DG. Statistical methods for assessing agreement between two methods of clinical measurement. *Lancet*. 1986; 307 – 10.
4. Bruton A, Conway JH, Holgate ST. 'Reliability: What is it and how is it measured?' *Physiotherapy*.2000; 86Suppl 2:94–99.
5. Caspersen CJ, Powell KE, Christenson GM. Physical activity, exercise, and physical fitness: definitions and distinctions for health-related research. *Public Health Rep*. 1985;100:126–31.
6. Cicchetti D, Bronen R, Spencer S, Haut S, Berg A, Oliver P, Tyrer P. Rating scales, scales of measurement, issues of reliability: resolving some critical issues for clinicians and researchers. *Journal of Nervous Mental Disease*. 2006;194:557–64.
7. Cicchetti DV. "Guidelines, criteria, and rules of thumb for evaluating normed and standardized assessment instruments in psychology". *Psychol Assess*. 1994;6(Suppl 4):284–90.
8. De Vet HCW, Terwee CB, Knol DL, Bouter LM. When to use agreement versus reliability measures. *J Clin Epidemiol*. 2006;59(Suppl 10):1033–9.
9. De Vet HC, Terwee CB, Mokkink LB, Knol DL. *Measurement in medicine: a practical guide*. Cambridge University Press; 2011.
10. Haley SM, Fragala-Pinkham MA. Interpreting change scores of tests and measures used in physical therapy. *Phys Ther*. 2006;86(Suppl 5):735–43.
11. Kottner J, Audige L, Brorson S, Donner A, Gajewski BJ, Hróbjartsson A, Roberts C, Shoukri M, Streiner DL. Guidelines for Reporting Reliability and Agreement Studies (GRRAS) were proposed. *Int J Nurs Stud*. 2011;48:661–71.
12. Lee P, Liu CH, Fan CW, Lu CP, Lu WS, Hsieh CL. The test–retest reliability and the minimal detectable change of the Purdue Pegboard Test in schizophrenia. *J Formos Med Assoc*. 2013;112(Suppl 6):332–7.
13. D'Olhaberriague L, Litvan I, Mitsias P, Mansbach HH. A reappraisal of reliability and validity studies in stroke. *Stroke*. 1996;27(Suppl 12):2331–6.
14. Ortega FB, Ruiz JR, Castillo MJ, Sjöström M. Physical fitness in childhood and adolescence: a powerful marker of health. Sweden. *International Journal of Obesity*. 2008;32(Suppl 1):1.
15. Portney LG, Watkins MP. *Foundations of Clinical Research: Applications to Practice*. 3rd ed. Upper Saddle River: Pearson/Prentice Hall; 2009.
16. Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability. *Psychological bulletin*. 1979;86(Suppl 2):420.
17. Smits-Engelsman BCM. *Performance and fitness battery for children: PERF-FIT, Manual*. Cape Town 2018.
18. Smits-Engelsman B, Cavalcante Neto JL, Draghi TTG, Rohr LA, Jelsma LD. Construct validity of the PERF-FIT, a test of motor skill-related fitness for children in low resource areas. *Research in Developmental Disabilities*.2020a.
19. Smits-Engelsman BCM, Bonney E, Neto JLC, Jelsma DL. Feasibility and content validity of the PERF-FIT test battery to assess movement skills, agility and power among children in low-resource settings. *BMC Public Health*. 2020b;20(Suppl 1):11–39. doi:10.1186/s12889-020-09236-w.
20. Strainer DL, Norman GR, Cairney J. *Health Measurement Scales: A practical guide to their development and use* 5 ed. Oxford: University Press; 2014.ISBN-13:9780199685219\$4 <https://doi.org/10.1093/med/9780199685219.001.0001> .
21. Stratford PW, Goldsmith CH. Use of the Standard Error as a Reliability Index of Interest: An Applied Example Using Elbow Flexor Strength Data. *Phys Ther*. 1997;77:745–50.
22. Valentini NC, Clark JE, Whitall J. Developmental co-ordination disorder in socially disadvantaged Brazilian children. *Child Care Health Dev*. 2015;41(Suppl 6):970–9.

23. Warburton DE, Gledhill N, Jamnik VK, Bredin SS, McKenzie DC, Stone J, Charlesworth S, Shephard RJ. Evidence-based risk assessment and recommendations for physical activity clearance: Consensus Document 2011. *Appl Physiol Nutr Metab.* 2011;36(Suppl 1):266-98.
24. Weir JP. Quantifying test-retest reliability using the intraclass correlation coefficient and the SEM. *The Journal of Strength Conditioning Research.* 2005;19(Suppl 1):231–40.
25. World Health Organization. International classification of functioning, disability and health—children & youth version. Geneva, Switzerland: 2007.

Figures

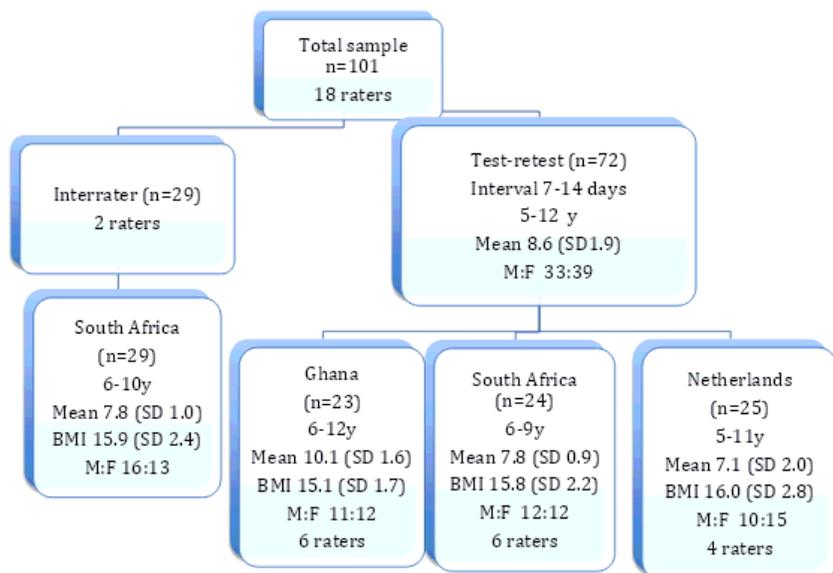


Figure 1

Flow chart for the total study. Country where the study took place and number of participants. Demographic participant information: Number of children per study, Age range, Mean Age and Standard Deviation (SD), Mean (SD) Body mass Index (BMI), Ratio Male (M): Female (F) participants and number of raters per study.

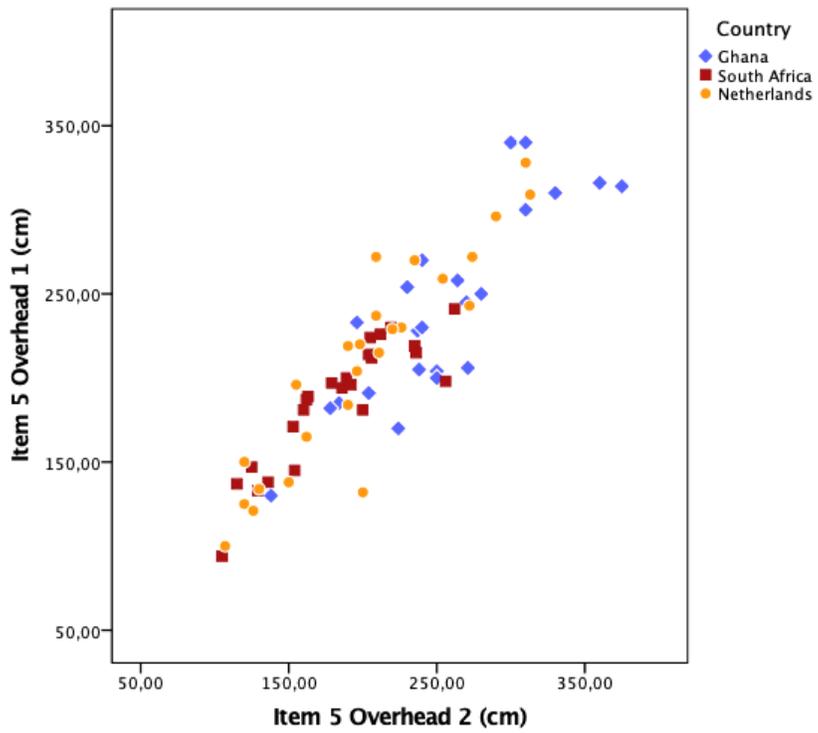


Figure 2

Scatterplot of the mean values (cm) obtained by the children for the Overhead throw at Time 1 (test) and Time 2 (retest) in the three countries.