

1 **PRE-PRINT**

2 **Inter-rater Reliability of the Test of Gross Motor Development – Third**
3 **Edition (TGMD-3) Following Raters’ Agreement on Measurement**
4 **Criteria**

5 **Abstract**

6 In this study, we calculated inter-rater reliability of the Test of Gross Motor
7 Development – Third Edition (TGMD-3) after raters reached a consensus
8 regarding measurement criteria. Three raters measured the fundamental movement
9 skills (FMS) of the same 25 children on the TGMD-3 at two different times, using
10 two different methods: (i) once when simply following the measurement criteria in
11 the TGMD-3 manual; and (ii) again, after a 9-month wash-out period, following
12 the raters’ consensus building for the measurement criteria for each skill. After
13 calculating and comparing the inter-rater reliability of these three raters across
14 these two rating times and methods, we found improved inter-rater reliability after
15 the raters’ consensus building discussions on ratings of both locomotor skills
16 (moderate-to-good reliability on 2 of 6 skills initially and at least moderate-to-
17 excellent on 4 of 6 skills following criteria consensus building) and ball skills
18 (moderate-to-good reliability on 1 of 7 skills initially and at least moderate-to-
19 excellent reliability on 4 of 7 skills following criteria consensus building). For
20 subtests scores and on overall test scores, raters achieved at least moderate-to-good
21 reliability on their second, post-consensus-building ratings. Based on this
22 improved reliability following consensus building, we recommend that researchers
23 include rater consensus-building before assessing children’s FMS or guiding
24 curriculum interventions in Physical Education from TGMD-3 data.

25
26 **Keywords:**

27 Assessment; child development; motor competence; gross motor skills; TGMD.
28

Introduction

The assessment of children's fundamental movement skills (FMS) has gained importance over recent decades (Bardid et al., 2019; Scheuer et al., 2019). FMS testing can be valuable for (a) identifying children with low levels of motor competence, permitting comparisons of motor proficiency levels across different populations, or guiding appropriate school interventions to promote children's healthy development (Scheuer et al., 2019; Tamplain et al., 2020).

A raft of FMS assessment tools is available for clinical, educational, and research purposes (Eddy et al., 2020). These tools can be broadly classified into (a) quantity/product-oriented test, those that offer quantifiable measurements of the product or outcome of children's movements (e.g., distance jumped); (b) quality/process-oriented tests, those that measure the quality of the movement process to determine whether children have yet attained some predefined behavioral criteria (e.g., judgments of arm movement quality during the jump); and (c) hybrid tests with scoring methods that combine both approaches (Bardid et al., 2019). Process-oriented tests can provide valuable qualitative information to guide teaching children *how to* accomplish a new motor movement (Stodden et al., 2008). However, scoring an observer's judgments of children's movements can be complex, and it may require raters to have extensive knowledge of FMS or specific skill training in their accurate assessment (Klingberg et al., 2019).

Among process-oriented assessment tools, the Test of Gross Motor Development-Third Edition (TGMD-3; Ulrich, 2019) (TGMD), and its predecessors TGMD and TGMD-2, is based on rater observations (Bardid et al., 2019) that are organized into two subscales – locomotor and ball skills. A recent systematic review suggested that the TGMD, in its various editions, is the most frequently used tool for measuring children's FMS

1 proficiency in educational, research, and clinical settings (Rey et al., 2020). A prerequisite
2 to the usage of any assessment tool is that it show adequate reliability and validity
3 (DeVellis, 2012; Streiner & Norman, 2008), and, despite touting these psychometric
4 strengths of the TGMD-3, Rey et al. (2020) also found in their review that inter-rater
5 reliability values were generally lower than those observed for intra-rater reliability,
6 probably due to rater difficulty achieving a consensus in the interpretation of scoring
7 criteria for some of the test's skill components (Barnett et al., 2014; Houwen et al., 2010).
8 Inter-rater differences relate to the rater's varied viewpoints, interpretations, and
9 assessment methods. The TGMD-3 test manual provides clear instructions to rate
10 whether a child meets certain performance criteria; but these judgments will always be
11 subject to each rater's discretion (Cano-Cappellacci et al., 2015), unless there has been
12 inter-rater consensus-building that is, ideally, widely disseminated and shared.

13 Barnett et al. (2014) examined the inter-rater reliability of the TGMD-2 object
14 control subtest by live observation and found some problematic, hard-to-identify
15 definitions of performance criteria that need to be clarified and discussed among raters in
16 a consensus building process. This TGMD weakness might imply a need for more
17 accurate TGMD measurement obtained by reducing the subjectivity bias of each rater
18 with criteria consensus prior to assessment with this tool. **Even after achieving**
19 **satisfactory inter-rater agreement, questions arise as to how this consensus is**
20 **sustained over time and when renewed training and consensus building may be**
21 **needed. This information is relevant, especially in studies in which measurements of**
22 **children's FMS span various periods, such as in intervention and longitudinal**
23 **research.** Yet, to the best of the authors' knowledge, no study has compared the inter-
24 rater reliability of TGMD after raters have reached agreement among themselves on

1 measurement criteria. Thus, our aim in this study was to determine inter-rater reliability
2 of the TGMD-3 after raters reached their own consensus on performance criteria.

3

4

Method

Participants

Twenty-five healthy primary school children participated in this study as TGMD-13 examinees (15 girls, 10 boys; $M = 9.16$, $SD = 1.31$ years). We obtained informed written consent from all the children's parents or guardians, and we obtained informed verbal assent from the participants. This study followed the Helsinki Convention's ethical principles and it was approved by the Ethical Committee of the Faculty of Education and Sport Sciences (University of xxxx).

Study Design

In a first stage, five raters used 13 video-recorded skills performed by 25 participants on the TGMD-13 to assess their motor competence (Carballo-Fazanes et al., 2021). Initially, each skill was explained, one by one, to the participants. Next, participants viewed a video, at normal speed and in slow motion, produced by the author of the TGMD, showing the correct execution of the skills (Ulrich & Webster, 2014). Participants then performed three trials of each skill: the first one was a practice trial permitting us to be sure they understood what they had to do, and the other two were video-recorded (camera Nikon D5300) so that the FMS could be measured by the raters later. We studied both intra-rater and inter-rater reliability by comparing novel versus expert raters' ratings when viewing the video continuously without pauses and at normal speed to viewing the video in slow-motion as many times as the rater needed (Carballo-Fazanes et al., 2021). Intra-rater reliability was higher than inter-rater reliability (both calculated throughout intraclass correlation coefficient (ICC)), regardless of the evaluator's experience or viewing mode, presumably due to the individual raters' separate viewpoints (Carballo-Fazanes et al., 2021).

1 In a second stage of this research, carried out nine months later, three out of the
2 five initial raters reached an agreement about assessing performance criteria of each skill
3 on the TGMD-3 following an inter-rater consensus-building process. After that, the three
4 raters measured the same 13 video-recorded skills of the same 25 participants included in
5 the first stage of the study.

6 ***Assessment Measure - Test of Gross Motor Development—Third Edition [TGMD-3]***

7 The TGMD-3 is a process-oriented test for assessing young children’s (aged 3-10
8 years) gross motor skill performance (Ulrich, 2019). It is organized into two subtests,
9 measuring locomotor and ball skills. The locomotor subtest measures skills that require
10 directional coordinated movements (run, gallop, hop, skip, horizontal jump, and slide).
11 The ball skills subtest measures skills related to intercepting and propelling objects (two-
12 hand strike, one-hand stationary dribble, overhand throw, kick, forehand strike, two-hand
13 catch, and underhand throw). Each skill includes 3-6 performance criteria with each one
14 scored as “0” or “1,” depending on the criterion’s absence or presence. Thus, a score is
15 obtained for each skill, for each subtest, and the sum of these item skills from both
16 subtests comprises the overall test score.

17 As noted, following a verbal description and video-based practical demonstration,
18 the participants performed three trials of each skill, with the first trial an unscored practice
19 trial to demonstrate the child’s understanding of the skill, and the other two rater-scored
20 trials, with “1” indicating correct performance and “0” indicating incorrect performance.
21 Scores ranged from 0 to 46 points for the locomotor subtest, and 0 to 54 for the ball skills
22 subtest, for an overall maximum score of 100.

23 ***Inter-Rater TGMD-3 Criteria Consensus Building***

24 Between their initial TGMD-3 ratings of participants and their second TGMD-3
25 ratings of the same participants after a 9-month wash-out period of their first ratings, our

1 three raters discussed TGMD-3 criteria together and reached a consensus for how to best
2 score the participants' performance on the TGDM-3. Their initial ratings (first stage)
3 relied upon the test manual guidelines for rating the 13 skills. In the second stage, **the**
4 raters **made individual revisions** of the **performance criteria for the** 13 skills and then
5 considered all their separate comments, skill by skill, and, through discussion together,
6 established both general and specific agreements (Table 1 and Table 2). They reached
7 general agreements on factors that could apply universally to all skills, and they reached
8 specific agreements that were applicable to only to a particular skill.

9 [Insert Tables 1 and 2 near here]

10 *Statistical Analyses*

11 We used the intraclass correlation coefficient (ICC) to assess inter-rater reliability.
12 Following Koo and Li (2016), we based ICC values and their 95% confidence intervals
13 on a single-measurement (type), consistency (definition), and 2-way random-effects
14 model. Values less than 0.50 indicated poor reliability, values between 0.50 and 0.75
15 indicated moderate reliability, values between 0.75 and 0.9 indicated good reliability, and
16 values greater than 0.90 indicated excellent reliability. Interpretation of ICC is based on
17 lower and upper bounds of the 95% confidence intervals.

18 We performed all analyses using the Statistical Package for the Social Sciences
19 (SPSS, version 23, IBM Corporation, Chicago, IL), and we set the statistical significance
20 level at $p < 0.05$.

21

22

Results

Inter-Rater Reliability of Locomotor Skills

The three raters' inter-rater reliabilities (ICC and 95% confident interval) of the participants' locomotor skills are shown in Table 3. Run and Hop skills had poor inter-rater reliability, both initially and after the raters' consensus-building process for performance criteria. For other locomotor skills, inter-rater reliability was higher after the consensus-building process in all pairwise rater comparisons. In the first stage, the inter-rater reliability calculated across the three raters' measurements was at least moderate in two skills (gallop and skip: moderate-to-good in both skills); however, on the second ratings following inter-rater consensus building, it was good-to-excellent for slide and skip and moderate-to-excellent for gallop and horizontal jump. Thus, scoring reliability for these three raters improved after the inter-rater consensus building

[Insert Table 3 near here]

Inter-Rater Reliability of Ball Skills

Inter-rater reliabilities (ICC and 95% confident interval) of ball skills are shown in Table 4. Inter-rater reliability increased following the raters' discussions about TGMD-3 performance criteria during consensus building. There was poor inter-rater reliability for kick and underhand throw skills on the initial ratings before consensus building (ICC < 0.5 in all comparisons), but ICC values associated with these two skills overcame the 0.50 threshold in all analyses after the raters' consensus building. A particularly pronounced improvement was evident in the case of underhand throw for which inter-rater reliability after consensus building was moderate-to-excellent reliability in all pairwise comparisons. Regarding inter-rater reliability calculated across the three raters' measurements, raters reached moderate-to-excellent reliability for three skills (two-hand strike, forehand strike and underhand throw) and good-to-excellent reliability

1 in 1 skill (overhand throw) after consensus building, while, initially, just 1 of the 7 ball
2 skills (two-hand strike) reached moderate reliability.

3 [Insert Table 4 near here]

4 ***Inter-Rater Reliability on TGMD-3 Subtest Scores and Overall Scores***

5 Inter-rater reliabilities (ICC and 95% confident interval) of subtest and overall
6 scores are shown in Table 5. Raters showed improvements when comparing initial ratings
7 to second ratings nine months later following consensus building. In the case of the
8 locomotor subtest, there was improved reliability in all analysis from initial ratings to
9 post-consensus building ratings (i.e., from “por to-good” reliability initially and
10 “moderate to good” on second ratings. Regarding the ball subtest, inter-rater reliability
11 was only worse between raters A and B (moderate-to-excellent) before the consensus
12 comparing with after consensus (good-to-excellent). Reliability of overall scores
13 remained the same for ratings made before and after the raters’ consensus building around
14 performance criteria.

15 [Insert Table 5 near here]

16

17

Discussion

In this study, we aimed to determine inter-rater reliability among raters of children's TGMD-3 motor skills before and after raters' consensus-building about performance skills' performance criteria. Inter-rater reliability improved after raters' consensus building on separate TGMD-3 locomotor skills and ball skills, and on TGMD-3 locomotor subtest scores and ball skills subtest scores.

More specifically, for locomotor skills, raters reached at least moderate inter-rater reliability when rating 4 out of the 6 skills after their consensus building, compared to reaching this level of inter-rater reliability for only 2 out of 6 skills on their initial ratings. After rater consensus-building, the most improved ratings were on the locomotor skills of slide and horizontal jump skills, for which inter-rater reliability improved from poor-to-good (slide) and poor-to-moderate (horizontal jump) initially to good-to-excellent (slide) and moderate-to-excellent (horizontal jump) on the second ratings. This significant improvement might demonstrate the importance, and value of having raters collectively objectify the assessment criteria. The locomotor skills with the lowest reliability were run and hop in both stages, and other research efforts have also documented lower inter-rater reliability for judging these skills (Carballo-Fazanes et al., 2021; Maeng et al., 2017; Valentini et al., 2017), possibly due to the extra subjectivity or complexity in their measurement criteria. For instance, one criterion in the run is "*narrow foot placement landing on heel or toes (not flat-footed)*", which is difficult to perceive in each stride in the live-assessments or even in video-assessments, since the camera should be far enough away to record 20 meters of running. Also, different prior studies have reported that the hop skill has the lowest inter-rater reliability on the TGMD-3 (Carballo-Fazanes et al., 2021; Rintala et al., 2017; Valentini et al., 2017). However, Maeng et al. (2017) found inter-rater reliability highest on the hop skill. On this skill, two performance criteria might

1 be particularly subjective: “*Non-hopping leg swings forward in pendular fashion to*
2 *produce force*” and “*Arms flex and swing forward to produce force*”. Especially the part
3 of these criteria describing “... *to produce force*” requires raters to differentiate between
4 swinging as a natural movement in balance and swinging specifically to produce force.

5 Previous research has also demonstrated assessment complexity for ball skills
6 (Barnett et al., 2014) and, as a result, ball skills have shown poorer inter-rater reliability
7 than locomotor skills (Carballo-Fazanes et al., 2021). However, in our study, even on
8 these more complex skills, raters improved their inter-rater reliability from their initial
9 ratings to after the consensus building process. In this regard, agreement between the
10 three raters were at least moderate just in 1 ball skill before the consensus building
11 process, reaching at least a moderate-to-excellent inter-rater reliability in the second stage
12 of the study.

13 In the current study, the most problematic ball skills for understanding inter-rater
14 reliability on the TGMD-3 were the one-hand stationary dribble, two hands catch, and
15 kick. Based on the existing literature, no agreement on the quality of inter-rater reliability
16 can be established for these skills. For instance, Barnett et al. (2014) obtained lower
17 reliability values in catch skill while the rater agreement for the kick was excellent.
18 Consistent with our findings, other studies (Carballo-Fazanes et al., 2021; Y. Kim et al.,
19 2012; Maeng et al., 2017; Rintala et al., 2017) found poor-to-moderate reliability for
20 kicking. Again, we consider that this could be due to the subjectivity of the criteria for
21 this “kill: “*Non-kicking foot placed close to the ball*”. What is considered “close to the
22 ball?” In our study, some raters only considered children “close to the ball” only if their
23 foot was just next to the ball, while others allowed some distance between the ball and
24 the foot. These individualistic interpretations might contribute to rating variance and
25 fluctuations in inter-rater reliability indices.

1 The same discordance applies to one-hand stationary dribble, for which our results
2 indicated poor inter-reliability in both stages. Yet, the rater agreement was excellent in
3 Maeng et al.'s (2017) research. On this skill, our raters expressed the greatest difficulty
4 with the following criteria "*Pushes the ball with fingertips (not slapping at the ball)*" and
5 "*Contacts ball with one hand at about waist level*". Despite conducting video-
6 assessments and being able to stop and slow-motion the action, our raters found such
7 short time-lapse movements difficult to code. In this case, our lower inter-rater reliability
8 index even after raters had a chance to build a performance criteria consensus might come
9 have from an incorrect application of one of their general agreements: "*If the rater has*
10 *doubts in assessing a performance criterion, it will be scored as "1"*" This general
11 agreement was established considering that FMS is part of the child's gross motor
12 development in which the movement assessment of a specific performance criteria should
13 not be as exacting as in the assessment of a technical movement of, for example, a sports
14 performance.

15 Concerning the subtest scores, our results showed higher inter-rater reliability on
16 these holistic indices than for the individual skills. Most previous studies also obtained
17 good-to-excellent inter-rater reliability on locomotor and ball subtests scores and the
18 overall score (Estevan et al., 2017; S. Kim et al., 2014; Lopes et al., 2016; Mohammadi
19 et al., 2019; Simons et al., 2008; Valentini, 2012; Valentini et al., 2017; Wagner et al.,
20 2016), and this level of inter-rater reliability was higher than in our study. In any case,
21 reliability related to subtests scores, in our study, increased after consensus building from
22 levels demonstrated before consensus building, especially on locomotor subtest scores
23 that improved from poor-to-good to moderate-to-excellent inter-rater reliability.

24 Our different findings from previous studies may be because we used more
25 demanding descriptors for our ICC values (Koo & Li, 2016). For example, we only

1 considered ICC values above 0.75 as good and those above 0.90 as excellent reliability.
2 In contrast, from a qualitative perspective, other studies considered ICC values above 0.6
3 as good and above 0.75 as excellent (Capio et al., 2016; Estevan et al., 2017; Houwen et
4 al., 2010; Mohammadi et al., 2019; Rintala et al., 2017). In effect, if we had interpreted
5 the ICC the same as these earlier studies, our inter-rater reliability would have also been
6 good-to-excellent on subtests and overall TGMD-3 scores. In addition, although some
7 studies only used ICC values for interpreting reliability (Ayán et al., 2019; Aye et al.,
8 2017; Barnett et al., 2014; Estevan et al., 2017; Houwen et al., 2010; S. Kim et al., 2014;
9 Y. Kim et al., 2012; Maeng et al., 2017; Rintala et al., 2017; Valentini et al., 2017), we
10 followed recommendations from Koo and Li (2016) and used both lower and upper values
11 of the 95% confidence interval. This, together with the fact that various investigators
12 performed other statistical tests for assessing reliability, such as kappa (Lopes et al., 2016)
13 or Pearson's coefficient (Simons et al., 2008), might explain why inter-rater reliability
14 was lower in the present study despite our raters' improved agreement. Although previous
15 investigations stated the need for inter-rater consensus building before assessment
16 (Barnett et al., 2014; Cano-Cappellacci et al., 2015; Houwen et al., 2010; Maeng et al.,
17 2017), our study is the first to show the benefits of reaching a consensus between raters
18 about performance criteria before they conducted FMS measurements.

19 ***Limitations and Directions for Further Research***

20 Some limitations to this study should be noted. First, we used a small sample of
21 25 healthy children as examinees, and our results should be replicated with a larger and
22 more diverse sample (e.g., varying socioeconomic backgrounds and laterality
23 preferences) and with a greater number of raters of different background experiences (e.g.
24 experts and non-experts in FMS assessment, teachers, and health care professionals). This
25 would permit a more precise understanding of rater agreement for motor competence

1 measurements with process-oriented test batteries. Also, of possible relevance, our 25
2 child participants were assessed by the raters twice. Since Ulrich (2019) suggested that
3 raters may have introduced some memory bias after a 14-day wash-out period, we
4 introduced an interval of 9-months in the present study to manage this potential confound.

5 Regarding general improvements in methodology that are needed in this line of
6 research, the use of different statistical analyses across studies complicates comparing
7 reliability results, and there is value in standardizing this approach. Of still greater
8 importance, rather than a building rater consensus on performance criteria in separate
9 studies, there is a need for common agreement and common training across all users of
10 the TGMD-3 such that inter-rater reliability is not only improved within a research team
11 but across research teams and across clinical practitioners.

12

Conclusion

In this study, we showed that rater's agreement on performance criteria might improve inter-rater reliability in assessing children's TGMD-3 skills. We recommend that in both research and clinical uses of the TGMD-3, to guide appropriate curriculum interventions in Physical Education, care should be taken to engage in consensus building to improve rater agreement about complex performance judgments on this subjective qualitative measure. Ideally, there should be a sharing of improved rating criteria across studies to allow for a common agreement and specific, standardized rater training for all users of the TGMD-3.

1 **References**

- 2 Ayán, C., Cancela, J. M., Sánchez-Lastra, M. A., Carballo-Roales, A. I., Domínguez-
3 Meis, F., & Redondo-Gutiérrez, L. (2019). Reliability and validity of the TGMD-2
4 battery in a Spanish population. *Revista Iberoamericana de Diagnostico y*
5 *Evaluacion Psicologica*, 51(1), 21–33. <https://doi.org/10.21865/RIDEP50.1.02>
- 6 Aye, T., Oo, K. S., Khin, M. T., Kuramoto-Ahuja, T., & Maruyama, H. (2017).
7 Reliability of the test of gross motor development second edition (TGMD-2) for
8 Kindergarten children in Myanmar. *Journal of Physical Therapy Science*, 29(10),
9 1726–1731. <https://doi.org/10.1589/jpts.29.1726>
- 10 Bardid, F., Vannozzi, G., Logan, S. W., Hardy, L. L., & Barnett, L. M. (2019). A
11 hitchhiker’s guide to assessing young people’s motor competence: Deciding what
12 method to use. *Journal of Science and Medicine in Sport*, 22(3), 311–318.
13 <https://doi.org/10.1016/j.jsams.2018.08.007>
- 14 Barnett, L. M., Minto, C., Lander, N., & Hardy, L. L. (2014). Interrater reliability
15 assessment using the Test of Gross Motor Development-2. *Journal of Science and*
16 *Medicine in Sport*, 17(6), 667–670. <https://doi.org/10.1016/j.jsams.2013.09.013>
- 17 Cano-Cappellacci, M., Leyton, F. A., & Carreño, J. D. (2015). Content validity and
18 reliability of test of gross motor development in Chilean children. *Revista de Saude*
19 *Publica*, 49. <https://doi.org/10.1590/S0034-8910.2015049005724>
- 20 Capio, C. M., Eguia, K. F., & Simons, J. (2016). Test of gross motor development-2 for
21 Filipino children with intellectual disability: Validity and reliability. *Journal of*
22 *Sports Sciences*, 34(1), 10–17. <https://doi.org/10.1080/02640414.2015.1033643>
- 23 Carballo-Fazanes, A., Rey, E., Valentini, N. C., Rodríguez-Fernández, J. E., Varela-
24 casal, C., Rico-Díaz, J., Barcala-furelos, R., & Abelairas-Gómez, C. (2021). Intra-
25 Rater (Live vs . Video Assessment) and Inter-Rater (Expert vs . Novice)

1 Reliability of the Test of Gross Motor Development — Third Edition.
2 *International Journal of Environmental Research and Public Health*, 18, 1652.

3 DeVellis. R.F. (2012). *Scale Development: Theory and Application* (3rd ed.), Sage
4 Publications, Thousand Oaks, CA.

5 Eddy, L. H., Bingham, D. D., Crossley, K. L., Shahid, N. F., Ellingham-Khan, M.,
6 Otteslev, A., Figueredo, N. S., Mon-Williams, M., & Hill, L. J. B. (2020). The
7 validity and reliability of observational assessment tools available to measure
8 fundamental movement skills in school-age children: A systematic review. *PLoS*
9 *ONE*, 15(8), e0237919. <https://doi.org/10.1371/journal.pone.0237919>

10 Estevan, I., Molina-García, J., Queralt, A., Álvarez, O., Castillo, I., & Barnett, L. (2017).
11 Validity and reliability of the Spanish version of the Test of Gross Motor
12 Development-3. *Journal of Motor Learning and Development*, 5(1), 69–81.
13 <https://doi.org/10.1123/jmld.2016-0045>

14 Houwen, S., Hartman, E., Jonker, L., & Visscher, C. (2010). Reliability and validity of
15 the TGMD-2 in primary-school-age children with visual impairments. *Adapted*
16 *Physical Activity Quarterly*, 27(2), 143–159. <https://doi.org/10.1123/apaq.27.2.143>

17 Kim, S., Kim, M. J., Valentini, N. C., & Clark, J. E. (2014). Validity and reliability of
18 the TGMD-2 for South Korean children. *Journal of Motor Behavior*, 46(5), 351–
19 356. <https://doi.org/10.1080/00222895.2014.914886>

20 Kim, Y., Park, I., & Kang, M. (2012). Examining rater effects of the TGMD-2 on
21 children with intellectual disability. *Adapted Physical Activity Quarterly*, 29(4),
22 346–365. <https://doi.org/10.1123/apaq.29.4.346>

23 Klingberg, B., Schranz, N., Barnett, L. M., Booth, V., & Ferrar, K. (2019). The
24 feasibility of fundamental movement skill assessments for pre-school aged
25 children. *Journal of Sports Sciences*, 37(4), 378–386.

1 <https://doi.org/10.1080/02640414.2018.1504603>

2 Koo, T. K., & Li, M. Y. (2016). A Guideline of Selecting and Reporting Intraclass
3 Correlation Coefficients for Reliability Research. *Journal of Chiropractic*
4 *Medicine, 15*(2), 155–163. <https://doi.org/10.1016/j.jcm.2016.02.012>

5 Lopes, V. P., Saraiva, L., & Rodrigues, L. P. (2016). Reliability and construct validity
6 of the test of gross motor development-2 in Portuguese children. *International*
7 *Journal of Sport and Exercise Psychology, 16*(3), 250–260.
8 <https://doi.org/10.1080/1612197X.2016.1226923>

9 Maeng, H., Webster, E. K., Pitchford, E. A., & Ulrich, D. A. (2017). Inter- and
10 intrarater reliabilities of the test of gross motor development—third edition among
11 experienced TGMD-2 raters. *Adapted Physical Activity Quarterly, 34*(4), 442–455.
12 <https://doi.org/10.1123/apaq.2016-0026>

13 Mohammadi, F., Bahram, A., Khalaji, H., Ulrich, D. A., & Ghadiri, F. (2019).
14 Evaluation of the psychometric properties of the Persian version of the test of
15 Gross Motor Development-3rd edition. *Journal of Motor Learning and*
16 *Development, 7*(1), 106–121. <https://doi.org/10.1123/JMLD.2017-0045>

17 Rey, E., Carballo-Fazanes, A., Varela-Casal, C., & Abelairas-Gómez, C. (2020).
18 Reliability of the test of gross motor development: A systematic review. *PLoS*
19 *ONE, 15*(7 July), 1–26. <https://doi.org/10.1371/journal.pone.0236070>

20 Rintala, P. O., Sääkslahti, A. K., & Iivonen, S. (2017). Reliability assessment of scores
21 from video-recorded TGMD-3 performances. *Journal of Motor Learning and*
22 *Development, 5*(1), 59–68. <https://doi.org/10.1123/jmld.2016-0007>

23 Scheuer, C., Herrmann, C., & Bund, A. (2019). Motor tests for primary school aged
24 children: A systematic review. *Journal of Sports Sciences, 37*(10), 1097–1112.
25 <https://doi.org/10.1080/02640414.2018.1544535>

- 1 Simons, J., Daly, D., Theodorou, F., Caron, C., Simons, J., & Andoniadou, E. (2008).
2 Validity and reliability of the TGMD-2 in 7-10-year-old Flemish children with
3 intellectual disability. *Adapted Physical Activity Quarterly*, 25(1), 71–82.
4 <https://doi.org/10.1123/apaq.25.1.71>
- 5 Stodden, D. F., Goodway, J. D., Langendorfer, S. J., Roberton, M. A., Rudisill, M. E.,
6 Garcia, C., & Garcia, L. E. (2008). A Developmental Perspective on the Role of
7 Motor Skill Competence in Physical Activity: An Emergent Relationship. *Quest*,
8 60(2), 290–306. <https://doi.org/10.1080/00336297.2008.10483582>
- 9 Streiner, D. L., & Norman, G.R. (2008). *Health Measurement Scales: A practical guide*
10 *to their development and use* (4th ed.), Oxford University Press, Oxford.
- 11 Tamplain, P., Webster, E. K., Brian, A., & Valentini, N. C. (2020). Assessment of
12 Motor Development in Childhood: Contemporary Issues, Considerations, and
13 Future Directions, *Journal of Motor Learning and Development*, 8(2), 391-409.
14 <https://journals.humankinetics.com/view/journals/jmld/8/2/article-p391.xml>
- 15 Ulrich, D. A. (2000). *Test of Gross Motor Development - 2* (Second Ed). Pro-Ed
16 Publishers.
- 17 Ulrich, D. A. (2019). *Test of Gross Motor Development - 3* (Third Ed). Pro-Ed
18 Publishers.
- 19 **Ulrich, D. A., & Webster, E. K. (2014). TGMD-3 Administration.**
20 **https://www.youtube.com/watch?v=9WggHyZpXl0&ab_channel=KipWebster**
- 21 Valentini, N. C. (2012). Validity and Reliability of the TGMD-2 for Brazilian Children.
22 *Journal of Motor Behavior*, 44, 275–280.
- 23 Valentini, N. C., Zanell, L. W., & Webster, E. K. (2017). Test of Gross Motor
24 Development – Third Edition: Establishing Content and Construct Validity for
25 Brazilian Children. *Journal of Motor Learning and Development*, 5(1), 15–28.

1 <https://doi.org/10.1123/jmld.2016-0002>

2 Wagner, M. O., Webster, E. K., & Ulrich, D. A. (2016). Psychometric Properties of the
3 Test of Gross Motor Development 3rd Edition (German translation) – Results of a
4 Pilot-Study. *Journal of Motor Learning and Development*, 5(1), 29–44.

5 <https://dx.doi.org/10.1123/jmld.2016-0006>

6

- 1 Table 1. Consensus reached regarding locomotor skills.
- 2 Table 2. Consensus reached regarding ball skills.
- 3 Table 3. Inter-rater reliability* among raters. Locomotor skills.
- 4 Table 4. Inter-rater reliability* among raters. Ball skills.
- 5 Table 5. Inter-rater reliability* among raters. Subtests & Overall.
- 6