

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email info.bmjopen@bmj.com

BMJ Open

Psychometric properties of gross motor assessment tools for children: a systematic review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-021734
Article Type:	Research
Date Submitted by the Author:	15-Jan-2018
Complete List of Authors:	Griffiths, Alison; Monash University, Physiotherapy; The Royal Children's Hospital, Physiotherapy Toovey, Rachel; Murdoch Childrens Research Institute, Developmental Disability and Rehabilitation Research; University of Melbourne, Physiotherapy Morgan, Prue; Monash University, Physiotherapy Spittle, Alicia; University of Melbourne, Physiotherapy; Murdoch Childrens Research Institute, Victorian Infant Brain Studies
Keywords:	PAEDIATRICS, REHABILITATION MEDICINE, validity, gross motor assessment, reliability

SCHOLARONE™
Manuscripts

Peer Review Only

Psychometric properties of gross motor assessment tools for children: a systematic review

Alison Griffiths^{1,2,3}, Rachel Toovey^{3,4}, Prue E. Morgan¹, Alicia J. Spittle^{3,4}

1. Department of Physiotherapy, School of Primary and Allied Health Care, Monash University, Frankston, Victoria, Australia

2. Department of Physiotherapy, The Royal Children's Hospital, Parkville, Victoria, Australia

3. Murdoch Children's Research Institute, Parkville, Victoria, Australia

4. Department of Physiotherapy, The University of Melbourne, Parkville, Victoria, Australia

Corresponding Author:

Name: A/Prof Alicia J. Spittle

Address: Department of Physiotherapy, The University of Melbourne, Level 7 Alan Gilbert Building, 161 Barry Street, Parkville, Vic 3010, Australia

Email: aspittle@unimelb.edu.au

Funding: This study was part-funded by grants from the National Health and Medical Research Council Career Development Fellowship (AJS) 1053767 and Centre of Research Excellence in Newborn Medicine 1060733 (AJS and AG) and the Victorian Government's Operational Infrastructure Support Program.

Financial disclosure: The authors have declared that they have no financial relationships relevant to this article.

Conflict of interest: The authors have no conflict of interest.

Keywords: paediatrics, reliability, validity, rehabilitation medicine, gross motor assessment

Abstract

Objective:

Gross motor assessment tools have a critical role in identifying, diagnosing and evaluating motor difficulties in childhood. The objective of this review was to systematically evaluate the psychometric properties and clinical utility of gross motor assessment tools for children 2-12 years.

Method:

A systematic search of MEDLINE, Embase, CINAHL and AMED was performed. Methodological quality was assessed with the CONsensus-based Standards for the selection of health status Measurement INstruments (COSMIN) checklist and an outcome measures rating form was used to evaluate reliability, validity and clinical utility of assessment tools.

Results:

Seven assessment tools from 37 studies/manuals met the inclusion criteria: Bayley Scale of Infant and Toddler Development-III (Bayley-III), Bruininks-Oseretsky Test of Motor Proficiency-2 (BOT-2), Movement Assessment Battery for Children-2 (MABC-2), McCarron Assessment of Neuromuscular Development (MAND), Neurological Sensory Motor Developmental Assessment (NSMDA), Peabody Developmental Motor Scales-2 (PDMS-2) and Test of Gross Motor Development-2 (TGMD-2). Methodological quality varied from poor to excellent. Validity and internal consistency varied from fair to excellent (α 0.5-0.99). The Bayley-III, NSMDA and MABC-2 have evidence of predictive validity. Test re-test reliability is excellent in the BOT-2 (ICC=0.80-0.99), PDMS-2 (ICC=0.97), MABC-2 (ICC=0.83-0.96) and TGMD-2 (ICC=0.81-0.92). TGMD-2 has the highest interrater (ICC 0.88-0.93) and intrarater reliability (ICC=0.92-0.99).

Conclusions:

The majority of gross motor assessments for children have good-excellent validity. Test-retest reliability is highest in the BOT-2, MABC-2, PDMS-2 and TGMD-2. The Bayley-III has the best predictive validity at 2 years of age for later motor outcome. None of the

1
2
3 1 assessment tools demonstrate good evaluative validity. Further research on evaluative gross
4
5 2 motor assessment tools are urgently needed.
6
7
8

9 3 Strengths and limitations of this study

- 10
11 4 • This systematic review comprehensively assesses methodological quality of included
12 studies using the COSMIN checklist.
13
14 5
15 6 • Results of this systematic review can provide guidance to clinicians when choosing
16 gross motor assessment tools based on test psychometric properties and clinical
17 utility.
18
19 8
20 9 • Areas for future research are identified including improving the evidence of inter and
21 intrarater reliability and responsiveness to change as well as the ascertainment of
22 predictive validity over a longer period of time.
23
24 11
25 12 • Only articles or test manuals written in English were included.
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 Introduction

2 Motor function promotes cognitive and perceptual development in children and contributes
3 to their ability to participate in their home, school and community environments¹. Motor
4 impairment can negatively affect activity and participation levels of children², which may
5 lead to lower levels of physical activity, fitness and health into adulthood³. While severe
6 motor deficits are usually diagnosed before 2 years of age, mild motor deficits may not
7 become evident until children are in preschool and primary school environments where
8 they are exposed to increasingly complex tasks and compared to their peers³. Identification
9 of motor difficulties is an important step towards support and intervention for the child and
10 their family.

11 Healthcare professionals and researchers require standardised assessment tools to identify,
12 classify and diagnose motor problems in children⁴. Further, assessment tools are essential
13 to monitor the effects of intervention⁴. There is no gold standard of motor assessment for
14 children and the available tests vary in their ease of use and interpretability in clinical and
15 research settings, and whether they are norm or criterion referenced⁵. Criterion referenced
16 tests are designed to be scored as items or criteria are demonstrated; meaning that the
17 score is a reflection of a child's competence on the test items. Most available assessments
18 however, are norm referenced, meaning that a child's results are reported in relation to a
19 specific population⁴. The characteristics of the normed population should be taken into
20 consideration when interpreting test results as environmental and cultural differences have
21 been found to affect motor development⁶.

22 Health professionals should be aware of the validity and reliability of assessment tools to
23 assist in their instrument selection and interpretation of results. Validity refers to "The
24 degree to which [an instrument] is an adequate reflection of the construct to be measured"
25⁷. If an instrument does not have adequate construct or content validity then it may not be
26 assessing the skills that it purports to. Reliability refers to "the degree to which the
27 measurement is free from measurement error"⁷, which is significant when interpreting
28 results. If a child is assessed as being significantly delayed in their gross motor skills, the
29 reliability of that tool indicates the likelihood that a result is due to error.

1
2
3 1 A systematic review in 2010 by Slater⁸ evaluated performance-based gross motor tests for
4 2 children with developmental coordination disorder, however it did not include the second
5 3 and most recent version of the Movement Assessment Battery for Children 2 (MABC-2),
6 4 which is widely used. Brown and Lalor⁹ suggested that as a result of the changes to the
7 5 original Movement Assessment Battery for Children (MABC) in age range, age bands,
8 6 materials and tasks, that the MABC-2 requires independent reliability and validity
9 7 assessment. Over the past eight years there has been a significant increase in the number of
10 8 papers assessing the psychometric properties of motor assessment tools in children. A
11 9 systematic review of these and previous papers is warranted, in order to add to our
12 10 understanding of the psychometrics of standardised gross motor assessment tools.

11 11 The primary aim of this systematic review is to identify and evaluate the clinical utility and
12 12 psychometric properties of gross motor assessment tools appropriate for use in preschool
13 13 and school age children from 2-12 years. The secondary aim of this review is to identify the
14 14 methodological quality of the included studies and areas for further research.

15 Method

16 16 A comprehensive search strategy was completed in databases OVID Medline (1996 to May
17 17 2017), CINAHL plus (1937 to July 2017), Embase (1974 – May 2017) and AMED (1985 – July
18 18 2017). The search strategy used MeSH terms and text words for ('child' or 'paediatric') and
19 19 ('motor skills' or motor activity' or 'gross motor' or 'psychomotor' or 'developmental
20 20 coordination disorder') and ('questionnaires' or 'outcome assessment' or 'instrument' or
21 21 'task performance') and ('reliability' or 'validity' or 'psychometrics'). Reference lists of
22 22 included articles were also screened to identify any additional papers. If full texts were
23 23 unavailable or further information required regarding availability of manuals authors were
24 24 contacted.

25 25 Assessment tools were included if they were 1. Discriminative, predictive or evaluative of
26 26 gross motor skills, 2. Assessed \geq two gross motor (e.g. balance, jumping etc.) items, 3. Able
27 27 to extract a meaningful gross motor sub-score, 4. Applicable to children 2-12 years of age, 5.
28 28 Criterion or norm referenced test with a standardised assessment procedure and 6.
29 29 Instructional manuals are published or commercially available.

1
2
3 1 Articles describing use of the assessment tool were included if $\geq 90\%$ of the study
4 population were within 2-12 years of age, it was available in English and if validity and/or
5 reliability of the assessment tool was reported.
6
7

8
9 4 Assessment tools were excluded if they met any of the following criteria 1. Questionnaires
10 or screening tools, 2. Only applicable to children with a specific diagnosis (e.g. cerebral
11 palsy, Down's syndrome), 3. Test manuals not available in English and 4. The version of the
12 test has been superseded.
13
14
15

16
17 8 Titles and abstracts were screened by the first author. The remaining papers were obtained
18 in full text and reviewed by two authors (AG, RT or PM) with selection based on inclusion
19 and exclusion criteria. Papers and assessment tools were included after agreement by both
20 raters, with conflicting decisions discussed until a consensus was reached.
21
22
23

24
25 12 Methodological assessment of the papers was completed using the four-point scale of the
26 COnsensus-based Standards for the selection of health status Measurement INstruments
27 (COSMIN) checklist¹⁰. The COSMIN incorporates three quality domains: Validity, Reliability
28 and Responsiveness consisting of nine measurement properties: content, construct, cross
29 cultural and criterion validity, hypothesis testing, internal consistency, reliability,
30 measurement error and responsiveness⁷ (Supplementary Table 1).
31
32
33
34

35
36 18 The overall score for each measurement property on the COSMIN checklist is determined by
37 a 'worse score counts' approach¹⁰. Each property is rated as excellent, good, fair or poor
38 methodological quality based on descriptive criteria. Data extraction and assessment of
39 methodological quality was performed independently by two assessors (AG and RT). In the
40 case of any uncertainty a third reviewer (AS) performed a COSMIN assessment and
41 disagreement was resolved through discussion.
42
43
44
45

46
47 24 A data extraction form for each assessment tool was adapted from the CanChild Outcome
48 Measures Rating Form to collate information on clinical utility, validity, reliability and
49 responsiveness¹¹. Clinical utility includes the cost of manuals, kits, training requirements,
50 time to administer the assessment and the ease of scoring. All reported values for reliability
51 were collected with Intraclass Correlation Coefficients (ICC) directly compared.
52
53
54
55

Results

Figure 1 provides details of study selection. Eight assessment tools were identified for inclusion; Bayley Scale of Infant and Toddler Development III (Bayley-III), Bruininks-Oseretsky Test of Motor Proficiency 2 (BOT-2), Movement Assessment Battery for Children 2 (MABC-2), McCarron Assessment of Neuromuscular Development (MAND), Neurological Sensory Motor Developmental Assessment (NSMDA), Peabody Developmental Motor Scales 2 (PDMS-2), and Test of Gross Motor Development 2 (TGMD-2). The corresponding manuals were then added to the final yield resulting in 30 papers and 7 manuals. Nineteen assessment tools were excluded (Supplementary Table 2).

The majority of assessment tools identified in this review are discriminative and most lend themselves towards use in a research setting. All norm referenced tools are from western countries and each identified test covers a different age range as shown in Table 1.

Most of the tools assess at least two domains of function: gross motor and fine motor skills, although the Bayley III and the NSMDA assess six to seven different domains of development. The TGMD-2 is the only tool that only assesses gross motor skills.

There is some consistency of items included within the gross motor skill subsets between tests. Most include a locomotion task such as walking, running or stair climbing; an object control or manipulation task such as throwing or catching a ball; and a static or dynamic balance task such as standing on one leg or hopping. The PDMS-2 and the MAND also include strength assessments (the PDMS-2 only in some age groups).

The number of gross motor items for assessment vary both within and between the tools (Table 1). For example, the number of items tested in the Bayley-III and the PDMS-2 depends on the age and ability of the child. Several assessments report criteria for describing gross motor delay, although all test manuals warn against diagnosing delay based on a single assessment.

Table 1. Gross Motor Assessment Tool Characteristics

Assessment Tool	Domains Tested	Gross motor components tested	Age range	Diagnostic criteria	Primary purpose	Secondary purpose	Type of test	Normative sample (year)
Bayley-III ¹²	Gross motor, fine motor, cognitive, communication, social/emotional, adaptive	Static postures, dynamic movement, balance	1 mth – 3 yrs	Developmental delay: <25th centile or below 2SD. *	Discriminative	Predictive, Evaluative, Research tool	Norm	1700 children from the USA (2000)
BOT-2 ¹³	Gross motor, fine motor	Coordination, balance, running speed and agility, strength	4 – 21 yrs	*	Discriminative Evaluative	Research tool	Norm	1520 children from the USA (2005)
MABC-2 ¹⁴	Gross motor, fine motor, balance	Aiming and catching, static and dynamic balance	3 – 16 yrs	Traffic light system: Green = normal, amber = 'at risk' and red = definite motor impairment (<15%). *	Discriminative Evaluative	Intervention planning, Research tool	Norm	1172 children from United Kingdom (2006)
MAND ¹⁵	Gross and fine motor	Coordination, jumping, static and dynamic balance	3 yrs – 25 yrs	NDI 70-85 = mild 55-69 = moderate <55 = severe disability *	Evaluative	Research tool	Norm	2000 3-35 yrs from the USA (1970's)
NSMDA ¹⁶	Gross Motor, Fine Motor, Neurological, Postural Development, Infant Patterns of Movement, Sensory Motor. †	Sitting, kneeling, walking, balance, running, hopping, jumping, catching, motor planning	1 mth – 6 yrs	Total score 6-8 normal, 9-11 minimal, 12-14 mild, 15-19 moderate, 20-25 severe, >25 profound disability *	Evaluative Discriminative	Predictive, Research tool	Criterion	N/A
PDMS-2 ¹⁷	Gross motor, fine motor	Stationary (standing balance, sit-ups, push-ups), locomotion (walking, running, jumping, hopping, etc.), object manipulation (kick, throw, hit, catch)	Birth – 5 yrs	*	Discriminative Evaluative	Predictive, Research tool	Norm	2003 USA and Canada (1997-8)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

TGMD-2 ¹⁹	Gross Motor	Locomotion (run, gallop, hop leap, jump, slide) and Object control (batting, dribbling, catch, kick, throw, roll)	3 – 10 yrs	*	Discriminative Evaluative	Outcome measure, research tool, intervention planning	Norm	1208 USA children (1997-1998)
-----------------------------	-------------	---	------------	---	---------------------------	---	------	-------------------------------

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁹ NDI, Neurodevelopmental Index; SD, Standard Deviation; mth, month; yrs, years *, Advisable to use clinical reasoning; †, requires some manual handling; USA, United States of America

For peer review only

1
2
3 1 The PDMS-2 is notable for the inclusion of credit towards incomplete skills in the scoring system.
4 2 Most other tests award a point or credit towards a skill only if it is demonstrated to the full
5 3 satisfaction of the stated criteria (score of 0 or 1). The PDMS-2 however is scored 0-2 allowing for 1
6 4 mark to be allocated as a child progresses towards a skill without mastering it. The NSMDA marking
7 5 criteria is somewhat more complicated with a system of scores 1-4 with a symbol of + denoting
8 6 hyperactive response and – a hyporeactive response. The PDMS-2, MABC-2, BOT-2, MAND, TGMD-2
9 7 and Bayley-III all require raw scores to be converted to a standard (or scaled) score based on tables
10 8 supplied in the manuals. For the BOT-2 this is a multiple step process which can then be converted to
11 9 both sex-specific or combined standard scores and percentile ranks. A summary of assessment tool
12 10 characteristics can be found in Table 1.

11 Clinical Utility

12 The clinical utility of the assessment tools is summarised in Table 2, while scoring and administration
13 13 is detailed in Supplementary Table 3. The shortest administration time is for the TGMD-2 which has
14 14 10 items; whilst most manuals report 20-60 minutes is required to complete an assessment. These
15 15 times are not inclusive of equipment set up, pack up and scoring, which varies depending on the
16 16 amount of equipment and complexity of the scoring process. All assessments require the user to be
17 17 familiar with the test before administration and to possess a high level of understanding of child
18 18 movement and development. The MABC-2 and PDMS-2 are the only assessments that come with
19 19 supporting material to guide intervention post assessment (when the complete kit is purchased).

20 Methodological quality

21 All articles were assessed using the COSMIN checklist to determine methodological quality. Several
22 22 studies were marked down for failing to report missing data, having small sample sizes and for using
23 23 inappropriate statistical methods. A summary of the articles and corresponding COSMIN
24 24 methodology rating is provided in Table 3.

25 Validity

26 The content and construct validity of the included assessment tools are summarised in
27 27 Supplementary Table 4. Most assessments were developed by or with input from experts in the
28 28 field, with most also performing literature reviews. Bruininks and Bruininks¹³ performed
29 29 comprehensive surveys, pilot, tryout and standardisation studies before finalising the BOT-2,
30 30 providing the most comprehensively reported content validity.

1
2
3 1 Construct validity was confirmed with factor analysis (either exploratory or confirmatory) in most
4 2 assessment tools. The MABC-2 and the TGMD-2 have the most evidence for construct validity, with
5 3 the MABC-2 requiring some changes to remain valid in the Chinese and Dutch speaking populations
6 4 ^{20 21}. The BOT-2, MABC-2 and TGMD-2 all provide evidence of discriminant validity in particular age
7 5 or diagnosis groups. The NSMDA has minimal assessment of construct validity in children over 2
8 6 years. The Bayley-III, NSMDA and MABC-2 are the only assessments that provide evidence of
9 7 predictive validity (Suppl. Table 5). Concurrent validity between the MABC-2, PDMS-2 and BOT-2 is
10 8 moderate to high, whilst the TGMD-2 is only weakly correlated with the MABC-2 ⁵ (Suppl. Table 5).
11 9 The PDMS-2, TMGD-2 and NSMDA report correlations with other criteria such as paediatrician
12 10 diagnosis, physical fitness or psychomotor/intelligence tests.
13 11
14 12
15 13

Table 2. Clinical Utility of Gross Motor Assessment Tools

Assessment Tool	Time to administer (min)	Test Procedure	Target Examiner population	Training	Equipment/Manual
Bayley-III ¹²	30-90	Therapist administers in standardised order	Paediatric health professionals early childhood specialists	Formal training not required. DVD, webinars and workshops available	Comprehensive manual/kit: £1089 Test kit provides most equipment
BOT-2 ¹³	40-60	Therapist administered in standardised order	Paediatric health professionals early childhood specialists	Formal training not required	Comprehensive manual/kit: £961 Test kit provides most equipment
MABC-2 ¹⁴	20-40	Therapist administers items in standardised order. Some flexibility allowed.	Research psychologists, OT, PT, Paediatricians	Formal training not required.	Comprehensive manual/ kit: £1191 Test kit provides most equipment
MAND ¹⁵	15-20	Therapist administers items in standardised order.	Professionals e.g. education, neurology, OT, PT, psychology etc.	Formal training not required.	Manual and test kit: £1366 includes equipment
NSMDA ¹⁶	20-45	Observation followed by therapist administration of test items.	PT, OT	Formal training not required (but is available)	Comprehensive manual: £35. Equipment not included
PDMS-2 ¹⁷	45-60 (20-30 for GM only)	Standardised procedure.	Paediatric health professionals, PE teachers, early intervention specialists	Formal training not required	Comprehensive manual/kit: £553 Includes some but not all equipment required
TGMD-2 ¹⁹	15-20	Standardised procedure.	Teachers, health professionals (OT, PT, doctors)	Formal training not required	Kit includes manual and record form: £128. Equipment not included

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition¹²; BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition¹³; MABC-2, Movement Assessment Battery for Children 2nd edition¹⁴; MAND, McCarron Assessment of Neuromuscular Development¹⁵; NSMDA, Neurological Sensory Motor Developmental Assessment¹⁶; PDMS-2, Peabody Developmental Motor Scales 2nd edition¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition¹⁹; GM, Gross motor; OT, Occupational Therapy; PT, Physiotherapy; PE, Physical Education

Table 3. Methodological quality of included articles

Test	First author, Year	Country	Population (Age, Diagnosis)	Internal consistency	Reliability	Measurement error	Content validity	Structural validity	Hypothesis testing	Cross- cultural validity	Criterion validity	Responsive - ness
BAYLEY III	Bayley ¹²	USA	1-42 mths	Fair	Fair	Good	Excellent	Good	Good	-	Good	-
	Spittle, et al. ⁴	Australia	2,4 yrs, Ex prem	-	-	-	-	-	-	-	Good	-
	Visser, et al. ²³	Netherlands	2.2-10.8 yrs, GDD, L.I.	-	-	-	Excellent	Poor	-	-	-	-
BOT-2	Wuang and Su ²⁴	Taiwan	4-12 yrs ID	Excellent	Excellent	Excellent	-	-	-	-	-	Fair
	Wuang, et al. ²⁵	Taiwan	3-6 yrs ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Bruininks and Bruininks ¹³	USA	4-21 yrs	Good	Fair (interrater) Fair (test-retest)	Good	Excellent	Good	-	-	Good	-
MABC-2 (AB 1)	Ellinoudis, et al. ²⁶	Greece	3-5.5 yrs	Excellent	Good	-	-	-	-	-	-	-
	Hua, et al. ²⁰	China	3-6 yrs	Excellent	Good	-	Excellent	Excellent	-	Poor	Excellent	-
	Logan, et al. ⁵	USA	3-6 yrs	-	-	-	-	-	Fair	-	Fair	-
	Smits-Engelsman, et al. ²⁷	Belgium	3-4 yrs	Poor	Poor	Poor	-	-	-	-	-	-
	Holm, et al. ²⁸	Norway	7-9 yrs	-	Fair (interrater) Poor (intrarater)	Poor	-	-	-	-	-	-
MABC-2 (AB 2)	Kita, et al. ²⁹	Japan	7-10 yrs	Excellent	-	-	-	-	-	Poor	-	-
	Griffiths, et al. ³⁰	Australia	4-8 yrs	-	-	-	-	-	-	-	Good	-
MABC-2	Henderson, et al. ¹⁴	UK	3-16 yrs	-	Fair	Good	Excellent	-	-	-	-	-
	Niemeijer, et al. ²¹	Netherlands + Belgium	-	-	-	-	-	-	-	Poor	-	-
	Schulz, et al. ³¹	U.K	3-16 yrs	-	-	-	Excellent	Good	-	-	-	-
	Valentini, et al. ³²	Brazil	3-13 yrs	Fair	Fair	-	Fair	Poor	-	Poor	Poor	-

	Wuang, et al. ²⁵	Taiwan	3-6 yrs, ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Wuang, et al. ³³	Taiwan	6-12 yrs DCD	Poor	Fair	Good	-	-	-	-	-	Fair
MAND	Hands, et al. ³⁴	Australia	10-17 yrs	-	-	-	-	Excellent	-	-	-	-
	McCarron ¹⁵	USA	7yrs	-	-	-	Fair	Poor	-	-	Poor	-
NSMDA	Danks, et al. ³⁵	Australia	2 + 4 yrs ELBW	-	-	-	-	-	-	-	Fair	-
	MacDonald and Burns ³⁶	Australia	2 + 4 yrs CP	-	-	-	-	Fair	-	-	Poor	-
	Burns, et al. ³⁷	Australia	1-24 mths VLBW	Poor	-	-	Poor	-	-	-	-	-
	Burns, et al. ³⁸	Australia	1-mnths VLBW	-	-	-	-	Poor	-	-	Fair	-
PDMS-2	Hua, et al. ²⁰	China	3-6 yrs.	Excellent	Good	-	Excellent	Excellent	-	Poor	Excellent	-
	Wuang, et al. ²⁵	Taiwan	3-6 yrs ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Folio and Fewell ¹⁷	USA	0-71 mnths	Good	-	Poor	Excellent	Good	Good	-	Poor	-
TGMD-2	Barnett, et al. ³⁹	Australia	4-8 yrs	-	Fair	-	-	-	-	-	-	-
	Farrokhi, et al. ⁴⁰	Iran	3-11 yrs	Fair	Fair	-	Fair	Fair	-	-	-	-
	Houwen, et al. ⁴¹	Netherlands	6-12 yrs VI	Fair	Fair	-	-	Fair	-	-	-	-
	Kim, et al. ⁴²	Korea	8-12 yrs ID	-	Poor	-	-	-	-	-	-	-
	Kim, et al. ⁴³	Korea	5-6 yrs	Poor	Fair	-	-	Poor	-	-	Poor	-
	Logan, et al. ⁵	USA	3-6 yrs	-	-	-	-	-	Fair	-	Fair	-
	Rudd, et al. ⁴⁴	Australia	6-12 yrs	-	-	-	-	Good	-	-	-	-
	Simons, et al. ⁴⁵	Belgium	7-10 yrs ID	Good	Good (interrater) Poor (test-retest)	-	Excellent	Good	Good	-	-	-
	Valentini ⁴⁶	Brazil	3-10 yrs	Poor	Fair (test-retest) Good (intra,	-	Excellent	Good	-	Fair	Good	-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

interrater)											
Wong and Yin Cheung ⁴⁷	China	3-10 yrs	-	-	-	-	Fair	-	-	-	-
Ulrich ¹⁹	USA	3-10 yrs	Good	Fair (test-retest)	Fair	Poor	Good	-	-	Fair	-
Poor (interrater)											

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁹ Mths, Months; yrs, years; DCD, Developmental Coordination Disorder; VI, Vision Impairment; ID, Intellectual Disability; GDD, global developmental delay; L.I, Language Impairment; ELBW, Extremely Low Birth Weight; VLBW, Very Low Birth Weight; CP, Cerebral Palsy; prem, premature; USA, United States of America

For peer review only

1 Reliability

2 Internal consistency of assessments are summarised in supplementary table 6. The BOT-2's high
3 internal consistency is well supported, including for children with an intellectual disability^{25,48}. The
4 MABC-2 appears to have lower internal consistency than the BOT-2, which may be related to the
5 limited number of test items (eight) on the MABC-2. The highest values for internal consistency for
6 the MABC-2 were obtained in specific populations (Intellectual disability and developmental
7 coordination disorder) with poor to fair methodology only. Conversely the highest quality articles
8 reported the lowest values, although it should be noted that these assessed age band 1 (3-6 years)
9 only. Internal consistency is reported to be high for the PDMS-2, while the MAND does not currently
10 have published internal consistency data in this age group. The TGMD-2 is reported by two good
11 quality (and four poor to fair quality) articles to have excellent internal consistency, including for
12 children with vision impairment and intellectual disability

13 The reliability findings are summarised in Supplementary Table 6 and in Figures 2 and 3. Test-retest
14 reliability was excellent in the Bayley-III (Supplementary Table 6), BOT-2 and PDMS-2; and was good
15 to excellent in the MABC-2 and TGMD-2 (Figure 2). Intra-rater reliability was rarely investigated or
16 reported for most tools, with the TGMD-2 demonstrating better results than the MABC-2 (Figure 3).
17 Only the TGMD-2 and MABC-2 report inter-rater reliability values using an ICC (Figure 3)^{28,39}. Inter-
18 rater reliability is also supported in the BOT-2 with Pearson Correlation Coefficient and Kappa
19 respectively. The studies referred to in the test manuals for the TGMD-2, Bayley-III, BOT-2 and
20 MABC-2 all report reliability findings using Pearson's correlation, which is less ideal than an ICC or
21 weighted kappa for statistical analysis^{49,50}. Only studies reporting ICC's are visually represented in
22 Figures 2 (test-retest) and 3 (inter and intra-rater). The TGMD-2 test-retest reliability results from
23 Houwen, et al.⁴¹ were believed to contain an error as the reported ICC was outside of the reported
24 confidence intervals (ICC 0.92, 0.82-0.91). This data set was therefore excluded from Figure 2.

25 Responsiveness was reported for the Bayley-III, BOT-2, MABC-2 and PDMS-2 with minimal
26 detectable change (MDC) or a standard error of measurement (SEM)²⁵. There have been no studies
27 to date on the responsiveness of the TGMD-2, NSMDA or MAND.

28 Discussion

29 This review identified eight gross motor assessment tools appropriate for use in clinical or research
30 settings, each with their own strengths and limitations. Interestingly, only one of the eight

1 assessments measured gross motor skills in isolation. This is likely a reflection on current practice to
2 assess children's development as a whole, rather than assessing individual domains in isolation.

3 The current review adds to the literature by including a thorough methodological assessment using
4 the COSMIN checklist. Our findings are consistent with an earlier review by Slater, et al.⁸ who
5 reported that the psychometric properties of the TGMD-2 and the BOT-2 were robust in children
6 with developmental coordination disorder. The MABC-2 and the PDMS-2 were also identified as
7 well supported assessment tools in this review. All assessment tools were found to have merits and
8 limitations and should be chosen with consideration for their psychometric properties, clinical utility
9 and for the population and age group in question.

10 Clinicians and parents who need guidance to set realistic therapy goals and to understand future
11 intervention requirements benefit from understanding a test's predictive ability. The NSMDA and
12 the MABC-2 are the only tools that have demonstrated long term (≥ 4 years follow up) predictive
13 validity, while the Bayley-III has good predictive validity at 2 years for future movement difficulties
14 and for the diagnosis of cerebral palsy at 4 years. However, further research into the long-term
15 predictive validity of all included gross motor assessment tools is warranted.

16 While validity and reliability should guide selection of assessment tools, clinical utility must also be
17 taken into consideration. Most tests have ongoing costs associated with forms and equipment
18 replacement, which may be prohibitive to some users. The NSMDA requires the therapist to handle
19 the child for several items which should be considered in relation to manual handling policies of
20 institutions. Assessment burden for children and families should also be taken into consideration
21 when selecting an assessment tool. Younger children are more likely to be distracted and may not
22 understand test items as well, which may also increase assessment times²⁷.

23 When a new edition of an assessment tool is released resulting in a change in age groups, scoring
24 or tasks it is insufficient to rely on the psychometric assessments that were performed on the
25 original test. The MABC-2 manual provides justification for the inclusion of reliability and validity
26 assessment of the original MABC¹⁴, however, owing to the significant changes in age groups and
27 tasks between editions these were not included for the analysis of the MABC-2 in this review. Two
28 studies quoted in the MABC-2 manual to support the validity and reliability are both unpublished
29 works and as such are also unable to be included in this systematic review. This could indicate a
30 publication bias for the MABC-2.

1 As yet there is little evidence to support the use of these assessments as outcome measures. The
2 TGMD-2 was created in part to be used as an outcome measure, however there are no articles to
3 date investigating its responsiveness to change¹⁹. The inclusion in some of the articles of minimal
4 detectable change (MDC) and minimal clinically important difference (MCID) is valuable for
5 clinicians. The difference between the two values is also of importance, as a change in score does
6 not necessarily relate to a meaningful change for the child or their family. It should also be noted
7 that all of the included assessment tools measure impairment and activity limitations, but do not
8 specifically address the other elements of the International Classification of Functioning, Disability
9 and Health (ICF) domains of participation, personal factors and environment². Clinicians should
10 utilise appropriate assessments or questionnaires to ensure that these domains of health are also
11 addressed in line with World Health Organisation guidelines².

12 When considering a test's reliability all three elements of test error should be taken into account –
13 these can be described as time sampling (assessed with test-retest reliability), content sampling
14 (assessed as internal consistency), and inter-scorer difference (or interrater reliability)¹⁹. This is one
15 of the reasons that clinicians should consider repeating assessments and/or completing a second
16 alternative assessment. All assessments should be interpreted in conjunction with clinical reasoning
17 and observation. Included assessment tools are not intended to be diagnostic on their own; results
18 need to be combined with other assessments and expert opinion to arrive at a clinical diagnosis.

19 In this review lower methodological scores on the COSMIN can be attributed to inadequate
20 reporting statistical methods, small sample sizes and non-independent assessors. Further research
21 in this area should consider addressing these limitations in their study design to reduce potential
22 error.

23 The thorough methodological assessment of the included articles using the COSMIN checklist
24 should be seen as a strength of this paper, as should the range of assessment tools included in this
25 review. While it has previously been argued that the 'worst score counts' criteria in the COSMIN
26 creates a floor effect⁵¹, the COSMIN authors argue that only 'fatal flaws' contribute to an overall
27 score of poor¹⁰. There are few tools available to assess the psychometric properties of assessment
28 tools and arguably none so robustly validated as the COSMIN.

29 There are many appropriate gross motor assessment tools available for use in research and clinical
30 settings today. The available tools demonstrate adequate validity and reliability and as such the
31 authors do not believe that new assessment tools need to be developed for use. There is scope

1 however to improve the evidence of inter and intrarater reliability and predictive validity should be
2 ascertained over a longer period of time and with greater methodological rigour. Tools also need
3 clearer assessment of their responsiveness to change to assist clinicians and researchers with
4 outcome measure selection. Researchers should be mindful of the methods they use to assess
5 validity and reliability. Clarity of reporting, statistical methods and sample sizes should be carefully
6 considered to ensure the highest quality of evidence.

7 Conclusion

8 Currently available motor assessment tools have good to excellent content and construct validity.
9 The BOT-2, MABC-2, PDMS-2 and TGMD-2 are the most reliable assessments in this age group. The
10 Bayley-III has the best predictive validity at 2 years of age, and the NSMDA and the MABC-2 both
11 have good predictive validity at 4 years of age. There is scope for further research into the
12 predictive validity, reliability and responsiveness of gross motor assessment tools in preschool and
13 school aged children. In practice clinicians should choose assessments with consideration of their
14 psychometric properties in the context of the child that they are assessing.

16 Author Contributions

17 All individuals listed as authors meet the appropriate authorship criteria and have approved the
18 acknowledgement of their contributions. AG was responsible for the drafting of the paper and
19 liaising with the co-authors on findings and conclusions. RT contributed to the paper through
20 interpretation of data, completing methodological assessments and revising the content throughout
21 its development. A/Profs PEM and AJS both contributed to the paper through assisting with the
22 development of research design, interpretation of data and revising the content through its
23 development.

Figures

Figure 1. PRISMA flow diagram detailing study selection

Figure 2. Test re-test reliability of gross motor assessment tools

Figure 2 legend: BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition ¹³; MABC-2, Movement Assessment Battery for Children 2nd edition ¹⁴; PDMS-2, Peabody Developmental Motor Scales 2nd edition ¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition ¹⁹.

Figure 3. Inter and interrater reliability of gross motor assessment tools

Figure 3 legend: MABC-2, Movement Assessment Battery for Children 2nd edition ¹⁴; TGMD-II, Test of Gross Motor Development 2nd edition ¹⁹

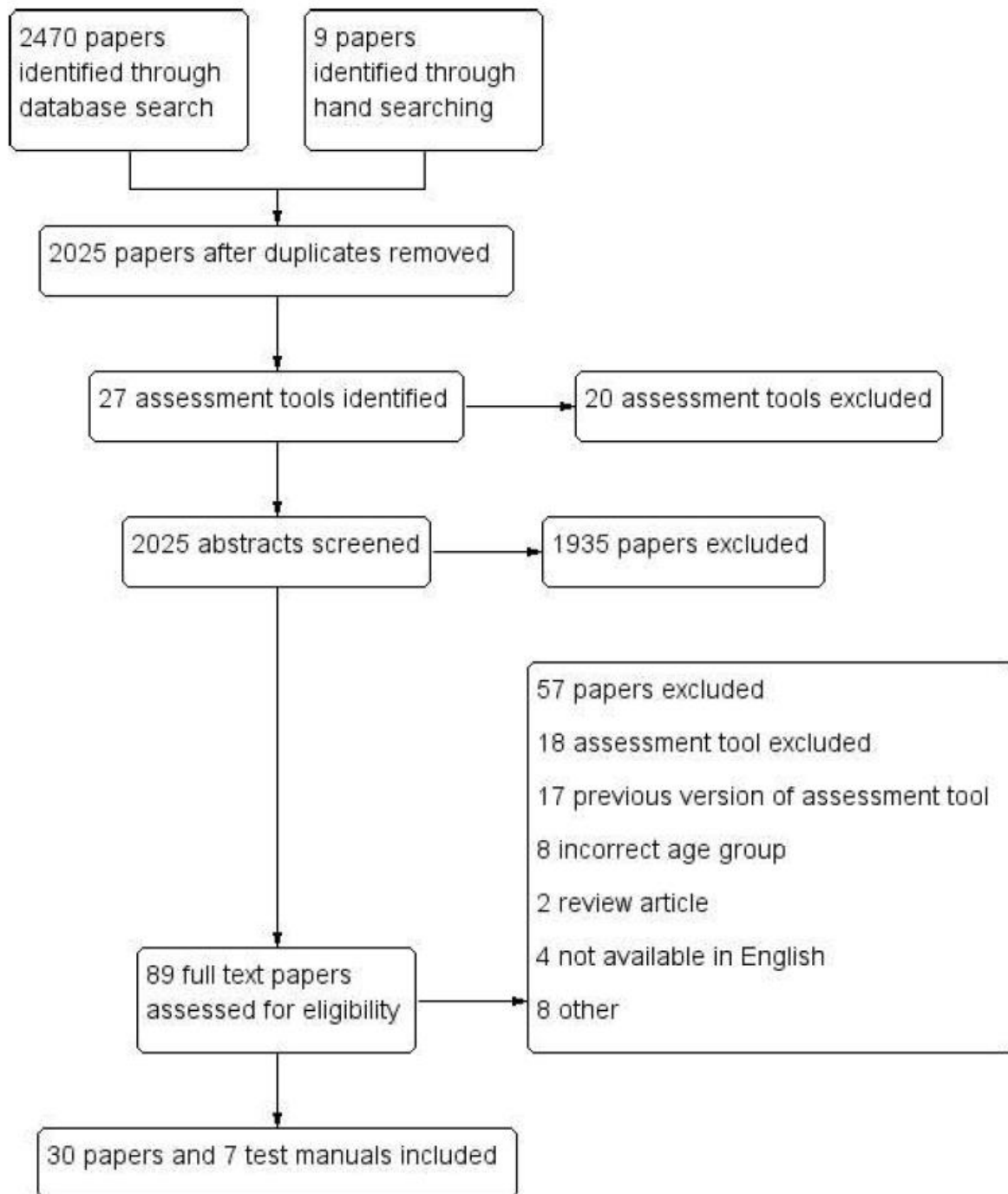
References

1. Piek JP, Baynam GB, Barrett NC. The relationship between fine and gross motor ability, self-perceptions and self-worth in children and adolescents. *Hum Mov Sci* 2006;25(1):65-75. doi: <http://dx.doi.org/10.1016/j.humov.2005.10.011>
2. World Health Organization. International Classification of Functioning, Disability and Health: ICF: World Health Organization 2001.
3. Magalhaes LC, Cardoso AA, Missiuna C. Activities and participation in children with developmental coordination disorder: a systematic review. *Res Dev Disabil* 2011;32(4):1309-16. doi: 10.1016/j.ridd.2011.01.029 [published Online First: 2011/02/19]
4. Spittle AJ, Spencer-Smith MM, Eeles AL, et al. Does the Bayley-III Motor Scale at 2 years predict motor outcome at 4 years in very preterm children? *Developmental Medicine And Child Neurology* 2013;55(5):448-52. doi: 10.1111/dmcn.12049
5. Logan SW, Robinson LE, Getchell N. The Comparison of Performances of Preschool Children on Two Motor Assessments. *Perceptual and Motor Skills* 2011;113(3):715-23.
6. Venetsanou F, Kambas A. Environmental factors affecting preschoolers' motor development. *Early Childhood Education Journal* 2010;37(4):319-27.
7. Mokkink LB, Terwee CB, Patrick DL, et al. The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. *J Clin Epidemiol* 2010;63(7):737-45. doi: 10.1016/j.jclinepi.2010.02.006 [published Online First: 2010/05/25]
8. Slater LM, Hillier SL, Civetta LR. The clinimetric properties of performance-based gross motor tests used for children with developmental coordination disorder: A systematic review. *Pediatr* 2010;22(2):170.
9. Brown T, Lalor A. The Movement Assessment Battery for Children—Second Edition (MABC-2): A Review and Critique. *Phys Occup Ther Pediatr* 2009;29(1):86-103.
10. Terwee CB, Mokkink LB, Knol DL, et al. Rating the methodological quality in systematic reviews of studies on measurement properties: a scoring system for the COSMIN checklist. *Qual Life Res* 2012;21(4):651-7. doi: 10.1007/s11136-011-9960-1 [published Online First: 2011/07/07]
11. Law M. Outcome measures rating form. Ontario, Canada: CanChild Centre for Disability Research, 2004.
12. Bayley N. Bayley Scales of Infant Development and Toddler Development: Technical Manual: The PsychCorp 2006.
13. Bruininks R, Bruininks B. Bruininks-Oseretsky Test of Motor Proficiency—2nd Edition (BOT-2): Manual. Circle Pines: MN: AGS Publishing 2005.
14. Henderson SE, Sugden DA, Barnett AL. Movement assessment battery for children-2: Movement ABC-2: Examiner's manual: Pearson 2007.

15. McCarron LT. MAND: McCarron assessment of neuromuscular development, fine and gross motor abilities: McCarron-Dial Systems, Incorporated 1997.
16. Burns YR. N.S.M.D.A Physiotherapy Assessment for Infants and Young Children Second Edition. Brisbane, Queensland: CopyRight Publishing Company 2014.
17. Folio M, Fewell R. Peabody Developmental Motor Scales. Examiner's Manual . 2nd Edition. Austin, Texas.: Pro-Ed. 2000.
18. Khan NZ, Muslima H, El Arifeen S, et al. Validation of a rapid neurodevelopmental assessment tool for 5 to 9 year-old children in Bangladesh. *J Pediatr* 2014;164(5):1165-70.e6. doi: <http://dx.doi.org/10.1016/j.jpeds.2013.12.037>
19. Ulrich DA. Test of gross motor development-2. *Austin: Prod-Ed* 2000
20. Hua J, Gu G, Meng W, et al. Age band 1 of the Movement Assessment Battery for Children-Second Edition: exploring its usefulness in mainland China. *Res Dev Disabil* 2013;34(2):801-8. doi: <http://dx.doi.org/10.1016/j.ridd.2012.10.012>
21. Niemeijer AS, van Waelvelde H, Smits-Engelsman BC. Crossing the North Sea seems to make DCD disappear: cross-validation of Movement Assessment Battery for Children-2 norms. *Hum Mov Sci* 2015;39:177-88. doi: 10.1016/j.humov.2014.11.004
22. Khan NZ, Muslima H, Shilpi AB, et al. Validation of rapid neurodevelopmental assessment for 2- to 5-year-old children in Bangladesh. *Pediatrics* 2013;131(2):e486-94. doi: <http://dx.doi.org/10.1542/peds.2011-2421>
23. Visser L, Ruiters SAJ, Van der Meulen BF, et al. Low verbal assessment with the Bayley-III. *Res Dev Disabil* 2015;36:230-43.
24. Wuang YP, Su CY. Reliability and responsiveness of the Bruininks-Oseretsky Test of Motor Proficiency-Second Edition in children with intellectual disability. *Res Dev Disabil* 2009;30(5):847-55. doi: <http://dx.doi.org/10.1016/j.ridd.2008.12.002>
25. Wuang YP, Su CY, Huang MH. Psychometric comparisons of three measures for assessing motor functions in preschoolers with intellectual disabilities. *J Intellect Disabil Res* 2012;56(6):567-78. doi: <http://dx.doi.org/10.1111/j.1365-2788.2011.01491.x>
26. Ellinoudis T, Evaggelina C, Kourtessis T, et al. Reliability and validity of age band 1 of the Movement Assessment Battery for Children--second edition. *Res Dev Disabil* 2011;32(3):1046-51. doi: <http://dx.doi.org/10.1016/j.ridd.2011.01.035>
27. Smits-Engelsman BCM, Niemeijer AS, van Waelvelde H. Is the Movement Assessment Battery for Children-2nd edition a reliable instrument to measure motor performance in 3 year old children? *Res Dev Disabil* 2011;32(4):1370-77. doi: <http://dx.doi.org/10.1016/j.ridd.2011.01.031>
28. Holm I, Tveter AT, Aulie VS, et al. High intra- and inter-rater chance variation of the movement assessment battery for children 2, ageband 2. *Res Dev Disabil* 2013;34(2):795-800. doi: <http://dx.doi.org/10.1016/j.ridd.2012.11.002>
29. Kita Y, Suzuki K, Hirata S, et al. Applicability of the Movement Assessment Battery for Children-Second Edition to Japanese children: A study of the Age Band 2. *Brain Dev* 2016;38(8):706-13.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
30. Griffiths A, Morgan P, Anderson PJ, et al. Predictive value of the Movement Assessment Battery for Children - Second Edition at 4 years, for motor impairment at 8 years in children born preterm. *Dev Med Child Neurol* 2017;59(5):490-96. doi: 10.1111/dmcn.13367 [published Online First: 2017/01/10]
31. Schulz J, Henderson SE, Sugden DA, et al. Structural validity of the Movement ABC-2 test: factor structure comparisons across three age groups. *Res Dev Disabil* 2011;32(4):1361-9. doi: <http://dx.doi.org/10.1016/j.ridd.2011.01.032>
32. Valentini NC, Ramalho MH, Oliveira MA. Movement assessment battery for children-2: translation, reliability, and validity for Brazilian children. *Res Dev Disabil* 2014;35(3):733-40. doi: <http://dx.doi.org/10.1016/j.ridd.2013.10.028>
33. Wuang YP, Su JH, Su CY. Reliability and responsiveness of the Movement Assessment Battery for Children-Second Edition Test in children with developmental coordination disorder. *Dev Med Child Neurol* 2012;54(2):160-5. doi: <http://dx.doi.org/10.1111/j.1469-8749.2011.04177.x>
34. Hands B, Larkin D, Rose E. The psychometric properties of the McCarron Assessment of Neuromuscular Development as a longitudinal measure with Australian youth. *Hum Mov Sci* 2013;32(3):485-97. doi: <http://dx.doi.org/10.1016/j.humov.2013.02.007>
35. Danks M, Maideen MF, Burns YR, et al. The long-term predictive validity of early motor development in "apparently normal" ELBW survivors. *Early Hum Dev* 2012;88(8):637-41.
36. MacDonald J, Burns Y. Performance on the NSMDA During the First and Second Year of Life to Predict Functional Ability at the Age Of 4 in Children with Cerebral Palsy. *Hong Kong Physiotherapy Journal* 2005;23(1):40-45. doi: 10.1016/S1013-7025(09)70058-2
37. Burns YR, Ensbeys RM, Norrie MA. The Neuro-sensory motor developmental assessment part 1: development and administration of the test. *Aust J Physiother* 1989;35(3):141-49.
38. Burns YR, Ensbeys RM, Norrie MA. The neuro-sensory motor developmental assessment part II: predictive and concurrent validity. *Aust J Physiother* 1989;35(3):151-57.
39. Barnett LM, Minto C, Lander N, et al. Interrater reliability assessment using the Test of Gross Motor Development-2. *Journal of Science and Medicine in Sport* 2014;17(6):667-70. doi: <http://dx.doi.org/10.1016/j.jsams.2013.09.013>
40. Farrokhi A, Zareh Zadeh M, Karimi Alvar L, et al. Reliability and validity of test of gross motor development-2 (Ulrich, 2000) among 3-10 aged children of Tehran City. *Journal of Physical Education and Sports Management* 2014;5(2):18-28. doi: 10.5897/JPEM12.003
41. Houwen S, Hartman E, Jonker L, et al. Reliability and validity of the TGMD-2 in primary-school-age children with visual impairments. *Adapt Phys Act Q* 2010;27(2):143-59.
42. Kim Y, Park I, Kang M. Examining rater effects of the TGMD-2 on children with intellectual disability. *Adapt Phys Act Q* 2012;29(4):346-65.
43. Kim CI, Han DW, Park IH. Reliability and validity of the test of gross motor development-II in Korean preschool children: applying AHP. *Res Dev Disabil* 2014;35(4):800-7. doi: <http://dx.doi.org/10.1016/j.ridd.2014.01.019>

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
44. Rudd J, Butson ML, Barnett L, et al. A holistic measurement model of movement competency in children. *J Sports Sci* 2016;34(5):477-85.
45. Simons J, Daly D, Theodorou F, et al. Validity and reliability of the TGMD-2 in 7-10-year-old Flemish children with intellectual disability. *Adapted physical activity quarterly : APAQ* 2008;25(1):71-82.
46. Valentini NC. Validity and reliability of the TGMD-2 for Brazilian children. *J Mot Behav* 2012;44(4):275-80. doi: <http://dx.doi.org/10.1080/00222895.2012.700967>
47. Wong KYA, Yin Cheung S. Confirmatory factor analysis of the Test of Gross Motor Development-2. *Measurement in Physical Education & Exercise Science* 2010;14(3):202-09. doi: 10.1080/10913671003726968
48. Wuang YP, Lin YH, Su CY. Rasch analysis of the Bruininks-Oseretsky Test of Motor Proficiency-Second Edition in intellectual disabilities. *Res Dev Disabil* 2009;30(6):1132-44. doi: <http://dx.doi.org/10.1016/j.ridd.2009.03.003>
49. Spittle AJ, Doyle LW, Boyd RN. A systematic review of the clinimetric properties of neuromotor assessments for preterm infants during the first year of life. *Dev Med Child Neurol* 2008;50(4):254-66. doi: 10.1111/j.1469-8749.2008.02025.x [published Online First: 2008/01/15]
50. McDowell I. Measuring health: a guide to rating scales and questionnaires: Oxford university press 2006.
51. Adair B, Said CM, Rodda J, et al. Psychometric properties of functional mobility tools in hereditary spastic paraplegia and other childhood neurological conditions. *Dev Med Child Neurol* 2012;54(7):596-605.



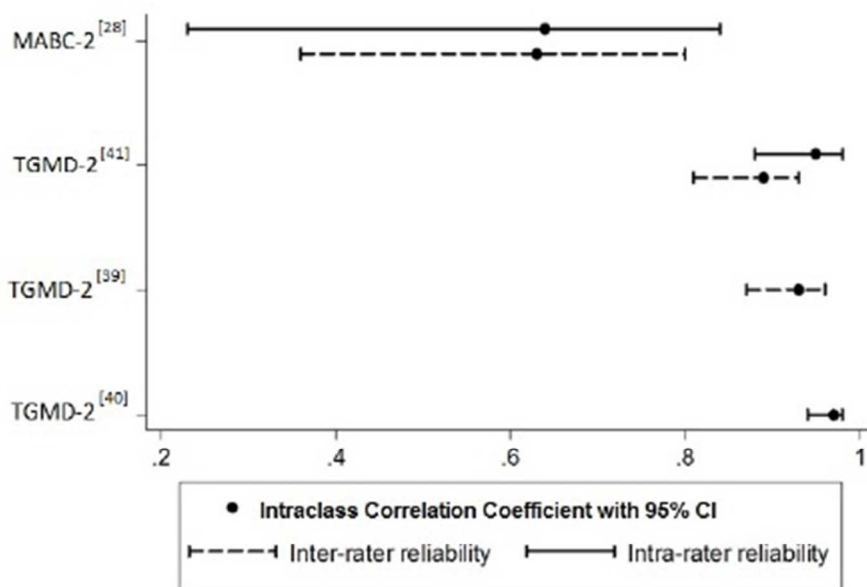


Figure 2. Test re-test reliability of gross motor assessment tools
 Figure 2 legend: BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition 13; MABC-2, Movement Assessment Battery for Children 2nd edition 14; PDMS-2, Peabody Developmental Motor Scales 2nd edition 17; TGMD-II, Test of Gross Motor Development 2nd edition 19.

41x28mm (300 x 300 DPI)

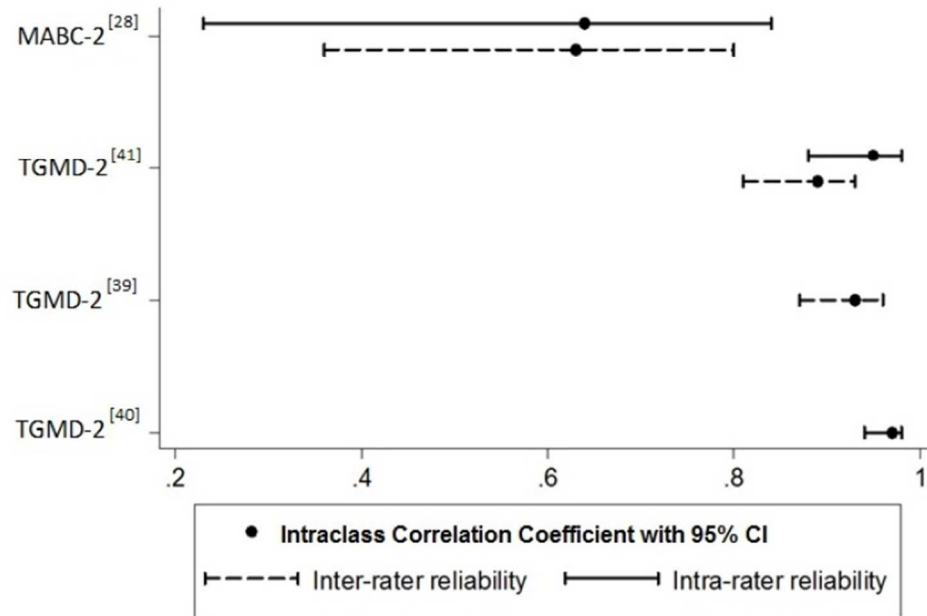


Figure 3. Inter and interrater reliability of gross motor assessment tools
 Figure 3 legend: MABC-2, Movement Assessment Battery for Children 2nd edition 14; TGMD-II, Test of Gross Motor Development 2nd edition 19

67x48mm (300 x 300 DPI)

Supplementary table 1: Definition of terms

		Definition	Example/explanation
Validity	Content	The degree to which an assessment tool's content measures the construct that it intends to measure [7]	Concerned with the relevance and comprehensiveness of the items included in the assessment tool
	Construct	Measures the degree to which the scores obtained from the test are an adequate reflection of the construct to be measured [7]	Examples include structural validity, hypothesis testing and cross-cultural validity
	Criterion	Assesses whether or not the test scores reflect a 'gold standard' assessment [7]	As there is no gold standard of assessment for gross motor function in children this is often assessed with correlations of scores obtained from two or three other frequently used tools.
Reliability		Refers to the consistency of a test score regardless of the time between assessments (test-retest) or the person administering (intra and inter-rater) [27]	Usually measured with intraclass correlation coefficient (ICC), but can be measured using Cohen's kappa coefficient. Percentage agreement and Pearson's correlation coefficient do not incorporate error into the calculations and as such is not a true measure of agreement [27]. Scores > 0.80 are considered excellent, 0.60-0.79 adequate and <0.59 poor [12]
	Internal consistency	The degree of interrelatedness of an assessment tool's items [7]	Usually measured using Cronbach's alpha (α) [7]. scores > 0.70 demonstrates high relationship, 0.5 to 0.69 a moderate relationship, 0.26 to 0.49 a low relationship and < 0.26

little relationship [27].

Measurement Error	Refers to the error obtained between measurements that cannot be attributed to the patients true change [7]	May be systematic or random error [7]
Responsiveness	An assessment tool's ability to detect change over time in the construct to be measured [7]	This is central to a tools capacity to be used as an outcome measure.

Supplementary table 2: Excluded Assessment Tools

Reason	Assessments
Manual not available in English	Maastricht's Motor Test (MMT) The Motor-Proficiency-Test for children between 4 and 6 years of age (MOT 4-6) Zuk Assessment Körperkoordinationstest für Kinder (KTK)
Cannot extract meaningful gross motor score	Early Intervention Developmental Profile (EIDP) Neurological Developmental Exam Preschooler Gross Motor Quality Scale (PGMQ) The Malawi Developmental Assessment Tool (MDAT) Dutch table tennis motor skills assessment
Screening Tool	Brief Assessment of Motor Function (BAMF) The Motor Performance Checklist Motor skill checklist (MSC)
Diagnosis specific/requires a diagnosis	Assessment Battery for the Atypical Handicapped Child (VAB) Video-based documentation and rating system of the motor behaviour of handicapped children
Only assesses one motor domain (e.g. gait)	Standardized Walking Obstacle Course (SWOC) Timed floor to stand test
Manual not published/commercially available	Rapid Neurodevelopmental Assessment (RNDA) Tufts Assessment of Motor Performance (TAMP) Zurich Neuromotor Assessment (ZNA)

Supplementary table 3: Scoring and administration of assessment tools

Assessment Tool	Scoring	Interpretation of scores	Other
Bayley-III [28]	Motor score - gross (varying items) and fine motor (varying items) subscales. Binary score with reverse/discontinue rules	Raw scores Composite scores Centile ranks Age equivalents Growth scores	Lends itself to multidisciplinary team testing.
BOT-2 [13]	Fine manual (15 items) manual coordination (12 items) body coordination (16 items) strength and agility (10 items) subscales. Scoring differs for subtests	Raw scores Age adjusted standard scores Composite scores Centile ranks Age equivalents Descriptive categories. Complex conversions	Administration Easel includes instructions, diagrams and photos of test procedure
MABC-2[24]	Manual dexterity (3 items), aiming & catching (2 items) and balance (3 items) subscales.	Raw scores component scores centile ranks total test score traffic light system. Simple conversion	Also Available: MABC-2 Checklist (screening tool) and intervention manual
MAND [29]	Fine motor (5 items) Gross motor (5 items)	Raw scores Scaled scores converted to an NDI. Factor scores. Complex conversions	Case studies included in manual for hyperactivity, encephalitis, mild head trauma, CP and muscular dystrophy
NSMDA [30]	Functional grade given for each subscale, which is combined to create an overall score.	Indicates: normal range, minimal dysfunction, mild problems, moderate, severe or profound disability	Sections for comment on strengths, behavioural state during testing, musculoskeletal system and recommendations.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

PDMS-2 [31]	GM: Stationary (30 items), locomotion (89 items), object manipulation (24 item). FM: grasping(26 items) , visual-motor integration (72 items)	Raw scores, Age equivalent, centile rank. Standard scores (subtests) Composite quotient. Complex conversions.	Motor activities program (intervention ideas)
TGMD-2 [25]	Locomotion (6 items) and Object Control (6 items). Separate male/female norms for object control subset	Raw scores, standard scores, percentile rank, age equivalent, Gross Motor Quotient. Simple conversion.	Simple to administer
Bayley-III, Bayley Scale of Infant and Toddler Development 3 rd edition [28]; BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2 nd edition [13]; MABC-2, Movement Assessment Battery for Children 2 nd edition [24]; MAND, McCarron Assessment of Neuromuscular Development [29]; NSMDA, Neurological Sensory Motor Developmental Assessment [30]; PDMS-2, Peabody Developmental Motor Scales 2 nd edition [31];; TGMD-II, Test of Gross Motor Development 2 nd edition [25]; GM, Gross Motor; FM, Fine Motor; NDI, Neurodevelopmental Index			

Supplementary table 4: Content and construct validity of assessment tools

Test	Content	Construct
BAYLEY III	Expert opinion for standard and low verbal version [28, 34]. Literature reviews. Gross motor score correlated with Motor component 0.70 [28]	Factor analysis. Difference in mean scores with pervasive developmental disorder, and specific language impairment [28]. H_i (gross motor subset) = 0.52-0.97 for children with language impairment and 0.82-0.99 in control group [34]
BOT-2	Focus groups, product survey, pilot, national tryout and standardisation studies, professional reviews[13]	Factor analysis, scores increase with age, discriminates between normal and children with DCD ($N=50$), high-functioning ASD ($N = 45$) and mild-moderate ID ($N = 66$) [13]
MABC-2	Expert Panel, Stakeholder feedback, Literature review [18] Expert panel - clarity (validity content index 71.8-93.9, Kappa 0.76-0.88) and pertinence (98.5-99.3 and kappa 0.83-0.92) $p<0.001$ [40]	Factor analysis, correlation coefficients [36] Subtest correlations 0.65-0.76 $p<0.001$. Discriminates between ASD and control group [18]. Structural equation modelling (for each age group) [39]. Expert panel - adequate face validity [40]. Significant difference between TD, DCD and at risk DCD scores ($\eta^2 = 0.63$) $p< 0.0001$ [40]. UK norms not appropriate to use with Dutch/Flemish children as under/over-estimate risk of motor impairment [15]. In Chinese population: CFA initially rejected. Acceptable fit achieved after 2 items removed [14]. Age band 2 shows good validity in Japanese population [37].
MAND	Based on neuropsychological theory. Several rounds of revision/trials of tasks during development [29]	Factor analysis [29] [42]. Scores increase with age, and discriminate between typically developing children and those with head trauma or neurological dysfunction as well as gender [29] [42]
NSMDA	Literature review. Developed by an experienced paediatric physiotherapist [45]	Factor analysis (up to 2 years of age) [45] [46]. Stability of test results over time (up to 2 years) [45] [46].
PDMS-2	Literature review. Created by experts in the field. Revised with feedback from therapists guided revision. Hierarchical sequence	Item response modelling. Factor analysis. Differential item functioning analysis. Scores correlated with age ($r=0.80-0.93$) [31]

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

of items [31]

TGMD-2	Expert Panel (3 PE teachers with post-grad qualifications) [25]. Translated version (Brazilian Portuguese) language clarity 0.96, pertinence >0.89. Experts CVI for clarity and pertinence were also strong- $\alpha = 0.93$ clarity and $\alpha = 0.91$ pertinence [52]	Exploratory and confirmatory factor analysis [50] [52] [25] [22] [53] [51] High and significant correlation of increasing age and increasing scores [47]. Age and disability differentiation [25] [51] Subtest correlation 0.41 [25] Gallop, running and leaping not well correlated with locomotion subscale. Object control significant & highly correlated [49]. ANOVA - significant age effect for object control [51] Moderate correlation between items and subset scores, and between subset scores and total score [51]
---------------	--	---

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;[28] BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;[13] MABC-2, Movement Assessment Battery for Children 2nd edition;[24] MAND, McCarron Assessment of Neuromuscular Development;[29] NSMDA, Neurological Sensory Motor Developmental Assessment;[30] PDMS-2, Peabody Developmental Motor Scales 2nd edition;[31] TGMD-II, Test of Gross Motor Development 2nd edition;[25]; H_i, scalability coefficient; CFA, Confirmatory Factor Analysis; TD, Typically Developing; ASD, Autism Spectrum Disorder, ID, Intellectual Disability; WPPSI, *Wechsler Preschool and Primary Scale of Intelligence*; WISC-R, Wechsler Preschool and Primary Scale of Intelligence-R; NDI, Neurodevelopmental Index; ANOVA, Analysis of Variance

Supplementary table 5: Criterion and predictive validity of assessment tools

Test	Concurrent/criterion	Predictive
BAYLEY III	Given but mean age <22 months. Not relevant to study population. [28]	Motor impairment at 4 years: Bayley III at 2 years <1SD = sensitivity 0.32-.037 specificity 0.97 <2SD sensitivity 0.18-0.21 specificity 1.00.

6

		CP at 4 years: Bayley III at 2 years <1SD sensitivity 0.83 specificity 0.94. <2SD sensitivity 0.67 specificity 1.0 [4]
BOT-2	MABC-2 $\rho = 0.92$ PDMS-2 $\rho = 0.88$ ($N = 38$) [17]. PDMS-2 Total motor composite $r = 0.77$ [13].	-
MABC-2	PDMS-2 $\rho = 0.631 - 0.84$ [17] [14]. TGMD-2 $\rho = 0.45$ [5]. TGMD-2 standard scores ($r = 0.3, p < 0.02$) [40]. BOT-2 $\rho = 0.90 - 0.92$ [17].	Classification groups (DCD, at risk and TD) remained same over time (6 months) $\chi^2 = 0.67 p = 0.72$ [40]. Predictive of motor impairment over 6-12 months ($N=41$) ICC 0.88 $p < 0.007$ [40]. Scores at 4 years predictive of motor impairment at 8 years in children born <30 weeks gestation (PPV 79, sensitivity 79%, specificity 93%) [38]
MAND	Gross motor subscore: Low-moderate correlation with manual dexterity (-0.46 to 0.35), reaction time (-0.31 to -0.58), intelligence measures (WISC-R, Metropolitan Achievement Test) (0.30-0.39) and visual motor test (-0.33 to 0.39) [29]	-
NSMDA	NSMDA at 2 years ($N = 148$) predictive of medical diagnosis $\chi^2 = 0.08 p = NS$ [46]	Motor outcome at 11-13 yrs. NSMDA at 2years - sensitivity 48.8%, specificity 82.4%, NSMDA at 4 years sensitivity 64.5%, and specificity 80%. PPV at 2 years 83% at 4 years 87% [43]. If classified 'severe' at 24 months - approximately 50% chance walking at 4 years (moderate = 80%, mild = 93% minimal = 100%) [44]
PDMS-2	MABC-2 $\rho = 0.63- 0.84$, [14, 17] MABC-2 gross motor composite $\rho = 0.743$ [14] BOT-2 $\rho = 0.88$ [17]. Mullen Scales of Early Learning GMQ = 0.86 FMQ = 0.80 [31]	-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

TGMD-2	MABC-2 total $r = 0.49$ $p < 0.01$ [5]. 'Teacher report' $r = 0.34-0.45$. physical fitness $r = -0.47 - 0.55$ [49]	-
	(N=41) Basic Motor Generalizations subtest of the CSSA $r = 0.63$. Locomotor 0.63 object control	
	0.41 [25]	

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition; [28] BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition; [13] MABC-2, Movement Assessment Battery for Children 2nd edition; [24] MAND, McCarron Assessment of Neuromuscular Development; [29] NSMDA, Neurological Sensory Motor Developmental Assessment; [30] PDMS-2, Peabody Developmental Motor Scales 2nd edition; [31] TGMD-II, Test of Gross Motor Development 2nd edition; [25] NS, Not Specified; SD, Standard Deviation; CP, Cerebral Palsy; TD, Typically Developing; ICC, Intraclass Correlation Coefficient; χ^2 , Chi Squared; NDI, Neurodevelopmental Index; CSSA, Comprehensive Scales of Student Abilities

Supplementary table 6: Reliability of assessment tools

Test	Internal Consistency	Test-Retest	Intra-rater	Inter-rater	Minimal detectable change	Minimal clinical important difference
BAYLEY III	GM $\alpha = 0.87-0.93$ MC: $\alpha 0.90-0.96$ (24-42 months) [28]	Gross Motor subtest (N=47) $r=0.79$ Motor component $r=0.80$ [28]	-	-	SEM Gross motor subtest 0.85-1.08. of Motor component = 3.00-4.74 (24-	-

8

42 months) [28]

BOT-2	(N = 100) $\alpha = 0.92$ [35] (N = 141) $\alpha = 0.86$ [17] 4-7 yrs (N= 620) $\alpha = 0.95$ 8-11 yrs (N= 450) $\alpha = 0.95$ [13]	(N = 100) ICC = 0.99 [35] (N = 141) ICC = 0.97 [17] 4-7 yrs (N = 43) $r = 0.81$ (8-12 yrs (N= 44) $r = 0.80$ [13]	-	Total motor composite 4-21 yrs (N = 47) $r = 0.98$ [13]	4.18 (sensitivity 55.10% specificity 72.55%) [35] 7.43 (sensitivity 42.49% specificity 65.72%) [17]	6.53 (sensitivity 48.98% specificity 76.47%) [35] 6.55 (sensitivity 49.99% specificity 58.78%) [17]
MABC-2 (AB 1)	(N = 60) M.D $\alpha = 0.51$, A&C $\alpha = 0.70$, Bal $\alpha = 0.66$ [36] (N = 1823) $\alpha = 0.502$ [14] (N=50) $\alpha = 0.81-0.87$ [23]	(N=60) ICC = 0.85 [36] Item ICC's 0.830-0.985 [14] ICC test-retest = 0.83 [23] Inter-rater test-retest ICC = 0.79 [23]	(N=28) $\kappa = 0.71$ [23]	Item ICC's range 0.892-0.998 [14] (N=22) $\kappa = 0.60$ [23]	(N=28) Intrarater MDC = 3.43 (N=22) Inter-tester MDC = 3.81 [23]	-
MABC-2 (AB 2)	Translated version (Japanese) (N=132) $\alpha = 0.602$ [37]	-	ICC = 0.64 [18]	ICC 0.63 [18]	Intra-rater SDC TTS: +/- 11.7 TSS +/- 3.3. Inter-rater SDC TTS +/-16.0 TSS +/- 3.8 [18]	-
MABC-2	Subscales $\alpha = 0.78$ (M.D = 0.77, BS = 0.52, Bal = 0.77) [40] $\alpha = 0.88$ [41] (N = 141) $\alpha = 0.88$ [17]	N=60 (all 3 age bands) $r=0.80$ [24] $r=0.74$ $p<0.0001$ (standard score). ICC standard score = 0.85 [40] ICC 0.96 [41] N = 141 ICC =0.96 [17]	ICC 0.88 [40]	ICC 0.96-0.99 [40]	SEM 1.34 (95%CI) = 3 [24] 1.83 (95%CI) [41] 1.83 (sensitivity 69.69% specificity 52.10%) [17]	1.39 (sensitivity 72.47% specificity 46.18%) [17, 41]

5	MAND	-	-	-	-	-	-
7	NSMDA	Cross correlation matrix Item scoring (12+24months) 0.73 $p < 0.001$, Functional grade (12+24months) 0.87 $p < 0.001$ [45]	-	-	-	-	-
13	PDMS-2	($N=141$) $\alpha=0.89$ [17] 24-35m $\alpha=0.97$, 36-47m $\alpha=0.95$, 48-59m $\alpha=0.97$, 60-71m $\alpha=0.98$. For subgroups† $\alpha=0.99$ [31]	$N=141$ ICC= 0.97 [17]	unable to extract data for ≥ 24 months [31]	unable to extract data for ≥ 24 months [31]	7.76 (sensitivity 60.65% specificity 74.13%) [17] SEM 24-59 months = 3, 60-71m = 2 [31]	8.39 (sensitivity 61.65% specificity 71.34%) [17]
19	TGMD-2	($N=1438$) $\alpha=0.80$ [47] $N=75$ Locomotor subset $\alpha=0.71$ object control $\alpha=0.72$ [22] $N=120$ $\alpha = 0.72$ [49] $N= 99$ $\alpha = 0.90$ [51] $N = 1208$ Cronbach's $\alpha = 0.91$ (gross motor quotient). Locomotor 0.85 and object control 0.88. Note SEM GMQ = 4-5 SEM subsets=1 [25]	$N=63$ ICC=0.81 95% CI [47] $N=23$ ICC=0.92 total 95% CI [22] $N=99$ $r=0.98$ [51] Locomotor test $r = 0.90$ $p < 0.0001$ object control test $r = 0.91$ $p < 0.001$ [52] $N = 75$ $r=0.96$ overall (3-5 yrs $r = 0.91$), 6-8 years $r = 0.95$), (9-10 years $r = 0.94$) [25]	$N=32$ ICC=0.97 95% CI [47] $N=25$ ICC=0.95 95% CI [22] ICC = 0.78 [48]	Obj ICC=0.93 [19] ($N=50$) ICC=0.89 [22] ICC=0.75 [48] $N=8$ $r=1.00$ [51] L.S ICC=0.88 Obj ICC=0.89 [52] $N = 30$ $r=0.98$ [25]	-	-

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;[28] BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;[13] MABC-2, Movement Assessment Battery for Children 2nd edition;[24] MAND, McCarron Assessment of Neuromuscular Development;[29] NSMDA, Neurological Sensory Motor Developmental Assessment;[30] PDMS-2, Peabody Developmental Motor Scales 2nd edition;[31] TGMD-II, Test of Gross Motor Development 2nd edition;[25] GM, Gross Motor Subset; MC, Motor Component; K, Kappa Coefficient; M.D, Manual Dexterity; BS, Ball Skills; BAL, Balance; A&C, Aiming and catching; SDC, Smallest Detectable Change; TTS, Total Test Score; TSS, Total Standard Score; †, gender, ethnicity, speech/language or physical disorder; Obj, Object Control Subset; L.S, Locomotion Subset



PRISMA 2009 Checklist

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

Section/topic	#	Checklist item	Reported on page #
TITLE			
Title	1	Identify the report as a systematic review, meta-analysis, or both.	1 Title page
ABSTRACT			
Structured summary	2	Provide a structured summary including, as applicable: background; objectives; data sources; study eligibility criteria, participants, and interventions; study appraisal and synthesis methods; results; limitations; conclusions and implications of key findings; systematic review registration number.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known.	4-5
Objectives	4	Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS).	5
METHODS			
Protocol and registration	5	Indicate if a review protocol exists, if and where it can be accessed (e.g., Web address), and, if available, provide registration information including registration number.	-
Eligibility criteria	6	Specify study characteristics (e.g., PICOS, length of follow-up) and report characteristics (e.g., years considered, language, publication status) used as criteria for eligibility, giving rationale.	5-6
Information sources	7	Describe all information sources (e.g., databases with dates of coverage, contact with study authors to identify additional studies) in the search and date last searched.	5
Search	8	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated.	5
Study selection	9	State the process for selecting studies (i.e., screening, eligibility, included in systematic review, and, if applicable, included in the meta-analysis).	5-6
Data collection process	10	Describe method of data extraction from reports (e.g., piloted forms, independently, in duplicate) and any processes for obtaining and confirming data from investigators.	6
Data items	11	List and define all variables for which data were sought (e.g., PICOS, funding sources) and any assumptions and simplifications made.	6
Risk of bias in individual studies	12	Describe methods used for assessing risk of bias of individual studies (including specification of whether this was done at the study or outcome level), and how this information is to be used in any data synthesis.	6
Summary measures	13	State the principal summary measures (e.g., risk ratio, difference in means).	7



PRISMA 2009 Checklist

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

Synthesis of results	14	Describe the methods of handling data and combining results of studies, if done, including measures of consistency (e.g., I^2) for each meta-analysis.	-
----------------------	----	---	---

Page 1 of 2

Section/topic	#	Checklist item	Reported on page #
Risk of bias across studies	15	Specify any assessment of risk of bias that may affect the cumulative evidence (e.g., publication bias, selective reporting within studies).	-
Additional analyses	16	Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified.	-
RESULTS			
Study selection	17	Give numbers of studies screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally with a flow diagram.	Figure 1 + page 7
Study characteristics	18	For each study, present characteristics for which data were extracted (e.g., study size, PICOS, follow-up period) and provide the citations.	Table 1 – page 8 + Suppl table 3
Risk of bias within studies	19	Present data on risk of bias of each study and, if available, any outcome level assessment (see item 12).	10 + Table 3 – page 13-14
Results of individual studies	20	For all outcomes considered (benefits or harms), present, for each study: (a) simple summary data for each intervention group (b) effect estimates and confidence intervals, ideally with a forest plot.	10-11, 16 + Table 2 – page 12 + Figures 2 – 3 + Suppl tables 4-6
Synthesis of results	21	Present results of each meta-analysis done, including confidence intervals and measures of consistency.	-
Risk of bias across studies	22	Present results of any assessment of risk of bias across studies (see Item 15).	-
Additional analysis	23	Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]).	-



PRISMA 2009 Checklist

DISCUSSION			
Summary of evidence	24	Summarize the main findings including the strength of evidence for each main outcome; consider their relevance to key groups (e.g., healthcare providers, users, and policy makers).	
Limitations	25	Discuss limitations at study and outcome level (e.g., risk of bias), and at review-level (e.g., incomplete retrieval of identified research, reporting bias).	16-18
Conclusions	26	Provide a general interpretation of the results in the context of other evidence, and implications for future research.	19
FUNDING			
Funding	27	Describe sources of funding for the systematic review and other support (e.g., supply of data); role of funders for the systematic review.	1

From: Moher D, Liberati A, Tetzlaff J, Altman DG, The PRISMA Group (2009). Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA Statement. PLoS Med 6(7): e1000097. doi:10.1371/journal.pmed1000097

For more information, visit: www.prisma-statement.org.

Page 2 of 2

BMJ Open

Psychometric properties of gross motor assessment tools for children: a systematic review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-021734.R1
Article Type:	Research
Date Submitted by the Author:	19-Jun-2018
Complete List of Authors:	Griffiths, Alison; Monash University, Physiotherapy; The Royal Children's Hospital, Physiotherapy Toovey, Rachel; Murdoch Childrens Research Institute, Developmental Disability and Rehabilitation Research; University of Melbourne, Physiotherapy Morgan, Prue; Monash University, Physiotherapy Spittle, Alicia; University of Melbourne, Physiotherapy; Murdoch Childrens Research Institute, Victorian Infant Brain Studies
Primary Subject Heading:	Paediatrics
Secondary Subject Heading:	Rehabilitation medicine
Keywords:	PAEDIATRICS, REHABILITATION MEDICINE, validity, gross motor assessment, reliability

SCHOLARONE™
Manuscripts

Only

Psychometric properties of gross motor assessment tools for children: a systematic review

Alison Griffiths^{1,2,3}, Rachel Toovey^{3,4}, Prue E. Morgan¹, Alicia J. Spittle^{3,4}

1. Department of Physiotherapy, School of Primary and Allied Health Care, Monash University, Frankston, Victoria, Australia

2. Department of Physiotherapy, The Royal Children's Hospital, Parkville, Victoria, Australia

3. Murdoch Children's Research Institute, Parkville, Victoria, Australia

4. Department of Physiotherapy, The University of Melbourne, Parkville, Victoria, Australia

Corresponding Author:

Name: A/Prof Alicia J. Spittle

Address: Department of Physiotherapy, The University of Melbourne, Level 7 Alan Gilbert Building, 161 Barry Street, Parkville, Vic 3010, Australia

Email: aspittle@unimelb.edu.au

Funding: This study was part-funded by grants from the National Health and Medical Research Council Career Development Fellowship (AJS) 1053767 and Centre of Research Excellence in Newborn Medicine 1060733 (AJS and AG) and the Victorian Government's Operational Infrastructure Support Program.

Financial disclosure: The authors have declared that they have no financial relationships relevant to this article.

Conflict of interest: The authors have no conflict of interest.

Keywords: paediatrics, reliability, validity, rehabilitation medicine, gross motor assessment

Abstract

Objective:

Gross motor assessment tools have a critical role in identifying, diagnosing and evaluating motor difficulties in childhood. The objective of this review was to systematically evaluate the psychometric properties and clinical utility of gross motor assessment tools for children 2-12 years.

Method:

A systematic search of MEDLINE, Embase, CINAHL and AMED was performed between May and July 2017. Methodological quality was assessed with the COnsensus-based Standards for the selection of health status Measurement INstruments (COSMIN) checklist and an outcome measures rating form was used to evaluate reliability, validity and clinical utility of assessment tools.

Results:

Seven assessment tools from 37 studies/manuals met the inclusion criteria: Bayley Scale of Infant and Toddler Development-III (Bayley-III), Bruininks-Oseretsky Test of Motor Proficiency-2 (BOT-2), Movement Assessment Battery for Children-2 (MABC-2), McCarron Assessment of Neuromuscular Development (MAND), Neurological Sensory Motor Developmental Assessment (NSMDA), Peabody Developmental Motor Scales-2 (PDMS-2) and Test of Gross Motor Development-2 (TGMD-2). Methodological quality varied from poor to excellent. Validity and internal consistency varied from fair to excellent (α 0.5-0.99). The Bayley-III, NSMDA and MABC-2 have evidence of predictive validity. Test re-test reliability is excellent in the BOT-2 (ICC=0.80-0.99), PDMS-2 (ICC=0.97), MABC-2 (ICC=0.83-0.96) and TGMD-2 (ICC=0.81-0.92). TGMD-2 has the highest interrater (ICC 0.88-0.93) and intrarater reliability (ICC=0.92-0.99).

Conclusions:

The majority of gross motor assessments for children have good-excellent validity. Test-retest reliability is highest in the BOT-2, MABC-2, PDMS-2 and TGMD-2. The Bayley-III has the best predictive validity at 2 years of age for later motor outcome. None of the

1
2
3 1 assessment tools demonstrate good evaluative validity. Further research on evaluative gross
4
5 2 motor assessment tools are urgently needed.
6
7
8

9 3 Strengths and limitations of this study 10

- 11 4 • This systematic review comprehensively assesses methodological quality of included
12 studies using the COSMIN checklist.
13 5
- 14 6 • Results of this systematic review can provide guidance to clinicians when choosing
15 gross motor assessment tools based on test psychometric properties and clinical
16 utility.
17 7
- 18 8 • Areas for future research are identified including improving the evidence of inter and
19 intrarater reliability and responsiveness to change as well as the ascertainment of
20 predictive validity over a longer period of time.
21 9
- 22 10 • Only articles or test manuals written in English were included.
23 11
- 24 12 • Only one reviewer screened titles and abstracts for inclusion
25 13
26 14

1 Introduction

2 Motor function promotes cognitive and perceptual development in children and contributes
3 to their ability to participate in their home, school and community environments¹. Motor
4 impairment can negatively affect activity and participation levels of children², which may
5 lead to lower levels of physical activity, fitness and health into adulthood³. While severe
6 motor deficits are usually diagnosed before 2 years of age, mild motor deficits may not
7 become evident until children are in preschool and primary school environments where
8 they are exposed to increasingly complex tasks and compared to their peers³. Identification
9 of motor difficulties is an important step towards support and intervention for the child and
10 their family.

11 Healthcare professionals and researchers require standardised assessment tools to identify,
12 classify and diagnose motor problems in children⁴. Further, assessment tools are essential
13 to monitor the effects of interventions⁴. There is no gold standard of motor assessment for
14 children and the available tests vary in their ease of use and interpretability in clinical and
15 research settings, and whether they are norm or criterion referenced⁵. Criterion referenced
16 tests are designed to be scored as items or criteria are demonstrated; meaning that the
17 score is a reflection of a child's competence on the test items. Most available assessments
18 however, are norm referenced, meaning that a child's results are reported in relation to a
19 specific population⁴. The characteristics of the normed population should be taken into
20 consideration when interpreting test results as environmental and cultural differences have
21 been found to affect motor development⁶.

22 Health professionals should be aware of the validity and reliability of assessment tools to
23 assist in their instrument selection and interpretation of results. Validity refers to "The
24 degree to which [an instrument] is an adequate reflection of the construct to be measured"
25⁷. If an instrument does not have adequate construct or content validity then it may not be
26 assessing the skills that it purports to. Reliability refers to "the degree to which the
27 measurement is free from measurement error"⁷, which is significant when interpreting
28 results. If a child is assessed as being significantly delayed in their gross motor skills, the
29 reliability of that tool indicates the likelihood that a result is due to error.

1
2
3 1 A systematic review in 2010 by Slater⁸ evaluated performance-based gross motor tests for
4 2 children with developmental coordination disorder, however it did not include the second
5 3 and most recent version of the Movement Assessment Battery for Children 2 (MABC-2),
6 4 which is widely used. Brown and Lalor⁹ suggested that as a result of the changes to the
7 5 original Movement Assessment Battery for Children (MABC) in age range, age bands,
8 6 materials and tasks, that the MABC-2 requires independent reliability and validity
9 7 assessment. Over the past eight years there has also been a significant increase in the
10 8 number of papers assessing the psychometric properties of motor assessment tools in
11 9 children. A systematic review of these and previous papers is warranted, in order to add to
12 10 our understanding of the psychometrics of standardised gross motor assessment tools.

13
14
15
16
17
18
19
20
21 11 The primary aim of this systematic review is to identify and evaluate the clinical utility and
22 12 psychometric properties of gross motor assessment tools appropriate for use in preschool
23 13 and school age children from 2-12 years by assessing the methodological quality of the
24 14 included studies. The secondary aim of this review is to identify any areas for further
25 15 research.

31 32 16 Method

33
34
35 17 A comprehensive search strategy was completed in databases OVID Medline (1996 to May
36 18 2017), CINAHL plus (1937 to July 2017), Embase (1974 – May 2017) and AMED (1985 – July
37 19 2017) (Supplementary tables 1-4). The search strategy used MeSH terms and text words for
38 20 ('child' or 'paediatric') and ('motor skills' or 'motor activity' or 'gross motor' or
39 21 'psychomotor' or 'developmental coordination disorder') and ('questionnaires' or 'outcome
40 22 assessment' or 'instrument' or 'task performance') and ('reliability' or 'validity' or
41 23 'psychometrics'). Reference lists of included articles were also screened to identify any
42 24 additional papers. If full texts were unavailable or further information required regarding
43 25 availability of manuals authors were contacted.

44
45
46
47
48
49
50
51 26 Assessment tools were included if they were 1. Discriminative, predictive or evaluative of
52 27 gross motor skills, 2. Assessed \geq two gross motor (e.g. balance, jumping etc.) items, 3. Able
53 28 to extract a meaningful gross motor sub-score, 4. Applicable to children 2-12 years of age, 5.

1
2
3 1 Criterion or norm referenced test with a standardised assessment procedure and 6.
4 2 Instructional manuals are published or commercially available.
5
6
7 3 Articles describing use of the assessment tool were included if; $\geq 90\%$ of the study
8 4 population were within 2-12 years of age, it was available in English and if validity and/or
9 5 reliability of the assessment tool was reported.
10
11
12
13 6 Assessment tools were excluded if they met any of the following criteria 1. Questionnaires
14 7 or screening tools, 2. Only applicable to children with a specific diagnosis (e.g. cerebral
15 8 palsy, Down's syndrome), 3. Test manuals not available in English and 4. The version of the
16 9 test has been superseded.
17
18
19
20
21 10 Titles and abstracts were screened by the first author with any studies that clearly did not
22 11 meet inclusion criteria excluded. The remaining papers were obtained in full text and
23 12 reviewed by two authors (AG, RT or PM) with selection based on inclusion and exclusion
24 13 criteria. Papers and assessment tools were included after agreement by both raters, with
25 14 conflicting decisions discussed until a consensus was reached.
26
27
28
29
30
31 15 Methodological assessment of the papers was completed using the four-point scale of the
32 16 COnsensus-based Standards for the selection of health status Measurement INstruments
33 17 (COSMIN) checklist¹⁰. The COSMIN incorporates three quality domains: Validity, Reliability
34 18 and Responsiveness consisting of seven measurement properties: content, construct and
35 19 criterion validity, internal consistency, reliability, measurement error and responsiveness⁷
36 20 (Supplementary Table 5). Cross-cultural validity, structural validity and hypothesis testing
37 21 are all considered to be a component of construct validity⁷. Whilst predictive validity is
38 22 considered to be a component of content validity, it is reported on separately in this paper
39 23 for interpretability of results⁷.
40
41
42
43
44
45
46 24 The overall score for each measurement property on the COSMIN checklist is determined by
47 25 a 'worse score counts' approach¹⁰. Each property is rated as excellent, good, fair or poor
48 26 methodological quality based on descriptive criteria. Data extraction and assessment of
49 27 methodological quality was performed independently by two assessors (AG and RT). In the
50 28 case of any uncertainty a third reviewer (AS) performed a COSMIN assessment and
51 29 disagreement was resolved through discussion.
52
53
54
55
56
57
58
59
60

1
2
3 1 A data extraction form for each assessment tool was adapted from the CanChild Outcome
4 2 Measures Rating Form to collate information on clinical utility, validity, reliability and
5 3 responsiveness¹¹. Items chosen to represent the clinical utility of the assessment tools were
6 4 the cost of manuals, kits, training requirements, time to administer the assessment and the
7 5 ease of scoring. All reported values for reliability were collected, however, only those papers
8 6 reporting intraclass Correlation Coefficient (ICC) were directly compared.

7 Results

8 Figure 1 provides details of study selection. Seven assessment tools were identified for
9 inclusion; Bayley Scale of Infant and Toddler Development III (Bayley-III), Bruininks-
10 Oseretsky Test of Motor Proficiency 2 (BOT-2), Movement Assessment Battery for Children 2
11 (MABC-2), McCarron Assessment of Neuromuscular Development (MAND), Neurological
12 Sensory Motor Developmental Assessment (NSMDA), Peabody Developmental Motor Scales
13 2 (PDMS-2), and Test of Gross Motor Development 2 (TGMD-2). The corresponding manuals
14 were then added to the final yield resulting in thirty papers and seven manuals. Twenty
15 assessment tools were excluded (Supplementary Table 6).

16 The majority of assessment tools identified in this review are discriminative and most lend
17 themselves towards use in a research setting. All norm referenced tools are from western
18 countries and each identified test covers a different age range as shown in Table 1.

19 The TGMD-2 is the only tool that assesses gross motor skills in isolation and that focusses on
20 quality of performance. The other gross motor assessments were either in conjunction with
21 assessment of fine motor and/or balance (MAND, MABC-2, BOT-2 and PDMS-2) or as a
22 component of a developmental assessment (NSMDA, Bayley-III).

23 Despite the variability in test structures, there is some consistency of items included within
24 the gross motor skill subsets between tests. Most include a locomotion task such as walking,
25 running or stair climbing; an object control or manipulation task such as throwing or
26 catching a ball; and a static or dynamic balance task such as standing on one leg or hopping.
27 The PDMS-2, BOT-2 and the MAND also include strength assessments (the PDMS-2 only in
28 some age groups).

1
2
3 1 The number of gross motor items for assessment vary both within and between the tools
4 (Table 1). For example, the number of items tested in the Bayley-III and the PDMS-2
5 2 depends on the age and ability of the child. Several assessments report criteria for
6 3 describing gross motor delay, although all test manuals warn against diagnosing delay based
7 4 on a single assessment.
8 5
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

Table 1. Gross Motor Assessment Tool Characteristics

Assessment Tool	Domains Tested	Gross motor components tested	Age range	Diagnostic criteria	Primary purpose	Secondary purpose	Type of test	Normative sample (year)
Bayley-III ¹²	Gross motor, fine motor, cognitive, communication, social/emotional, adaptive	Static postures, dynamic movement, balance	1 mth – 3 yrs	Developmental delay: <25th centile or below 2SD. *	Discriminative	Predictive, Evaluative, Research tool	Norm	1700 children from the USA (2000)
BOT-2 ¹³	Gross motor, fine motor	Coordination, balance, running speed and agility, strength	4 – 21 yrs	*	Discriminative Evaluative	Research tool	Norm	1520 children from the USA (2005)
MABC-2 ¹⁴	Gross motor, fine motor, balance	Aiming and catching, static and dynamic balance	3 – 16 yrs	Traffic light system: Green = normal, amber = 'at risk' and red = definite motor impairment (<15%). *	Discriminative Evaluative	Intervention planning, Research tool	Norm	1172 children from United Kingdom (2006)
MAND ¹⁵	Gross and fine motor	Coordination, jumping, static and dynamic balance	3 yrs – 25 yrs	NDI 70-85 = mild 55-69 = moderate <55 = severe disability *	Evaluative	Research tool	Norm	2000 3-35 yrs from the USA (1970's)
NSMDA ¹⁶	Gross Motor, Fine Motor, Neurological, Postural Development, Infant Patterns of Movement, Sensory Motor. †	Sitting, kneeling, walking, balance, running, hopping, jumping, catching, motor planning	1 mth – 6 yrs	Total score 6-8 normal, 9-11 minimal, 12-14 mild, 15-19 moderate, 20-25 severe, >25 profound disability *	Evaluative Discriminative	Predictive, Research tool	Criterion	N/A
PDMS-2 ¹⁷	Gross motor, fine motor	Stationary (standing balance, sit-ups, push-ups), locomotion (walking, running, jumping, hopping, etc.), object manipulation (kick, throw, hit, catch)	Birth – 5 yrs	*	Discriminative Evaluative	Predictive, Research tool	Norm	2003 USA and Canada (1997-8)

1									
2									
3									
4									
5	TGMD-2 ¹⁸	Gross Motor	Locomotion (run, gallop, hop, leap, jump, slide) and Object control (batting, dribbling, catch, kick, throw, roll)	3 – 10 yrs	*	Discriminative Evaluative	Outcome measure, research tool, intervention planning	Norm	1208 USA children (1997-1998)
6									
7									
8									
9									

10 Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for
 11 Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental
 12 Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ NDI, Neurodevelopmental Index; SD, Standard Deviation; mth, month; yrs, years *, Advisable to use
 13 clinical reasoning; †, requires some manual handling; USA, United States of America

1
2
3 1 The PDMS-2 is notable for the inclusion of credit towards incomplete skills in the scoring system.
4 2 Most other tests award a point or credit towards a skill only if it is demonstrated to the full
5 3 satisfaction of the stated criteria (score of 0 or 1). The PDMS-2 however is scored 0-2 allowing for 1
6 4 mark to be allocated as a child progresses towards a skill without mastering it. The TGMD-2 is also
7 5 notable for its marking system, in which points are awarded for the quality of the action performed,
8 6 instead of satisfactory completion of the task only. These actions include preparatory movements
9 7 prior to running and jumping, or arm position during movements. The NSMDA marking criteria is
10 8 somewhat more complicated with a system of scores 1-4 with a symbol of “+” denoting hyperactive
11 9 response and “-” a hyporeactive response. The PDMS-2, MABC-2, BOT-2, MAND, TGMD-2 and
12 10 Bayley-III all require raw scores to be converted to a standard (or scaled) score based on tables
13 11 supplied in the manuals. For the BOT-2 this is a multiple step process which can then be converted to
14 12 both sex-specific or combined standard scores and percentile ranks. A summary of assessment tool
15 13 characteristics can be found in Table 1.

14 Clinical Utility

15 The clinical utility of the assessment tools is summarised in Table 2, while scoring and administration
16 16 is detailed in Supplementary Table 7. The shortest administration time is 15-20 minutes for the
17 17 TGMD-2 and the MAND; whilst most manuals report 20-60 minutes is required to complete an
18 18 assessment. These times are not inclusive of equipment set up, pack up and scoring, which varies
19 19 depending on the amount of equipment and complexity of the scoring process. All assessments
20 20 require the user to be familiar with the test before administration and to possess a high level of
21 21 understanding of child movement and development. The MABC-2 and PDMS-2 are the only
22 22 assessments that come with supporting material to guide intervention post assessment (when the
23 23 complete kit is purchased).

24 Methodological quality

25 All articles were assessed using the COSMIN checklist to determine methodological quality. Several
26 26 studies were marked down for failing to report missing data, small sample sizes and for using
27 27 inappropriate statistical methods. A summary of the articles and corresponding COSMIN
28 28 methodology rating is provided in Table 3.

Validity

The content and construct validity of the included assessment tools are summarised in Table 4.

Most assessments were developed by or with input from experts in the field, with most also performing literature reviews. Bruininks and Bruininks¹³ performed comprehensive surveys, pilot, tryout and standardisation studies before finalising the BOT-2, providing the most comprehensively reported content validity.

Construct validity was confirmed with factor analysis (either exploratory or confirmatory) in most assessment tools. The TGMD-2 has the most evidence for construct validity with several papers performing confirmatory and exploratory factor analysis^{19 20 18 21 22 23}. The MABC-2, BOT-2, Bayley-III, MAND and PDMS-2 had factor analysis performed only in one paper. The MABC-2 was shown to require changes to remain valid in the Chinese and Dutch speaking populations^{24 25}. The BOT-2, MABC-2 and TGMD-2 all provide evidence of the ability to discriminate between particular age or diagnosis groups, which can be considered to support their content validity. The NSMDA has minimal assessment of construct validity in children over 2 years. The Bayley-III, NSMDA and MABC-2 are the only assessments that provide evidence of predictive validity (Table 5). Concurrent validity between the MABC-2, PDMS-2 and BOT-2 is moderate to high, whilst the TGMD-2 is only weakly correlated with the MABC-2⁵ (Table 5). The PDMS-2, TMGD-2 and NSMDA report correlations with other criteria such as paediatrician diagnosis, physical fitness or psychomotor/intelligence tests.

Table 2. Clinical Utility of Gross Motor Assessment Tools

Assessment Tool	Time to administer (min)	Test Procedure	Target Examiner population	Training	Equipment/Manual
Bayley-III ¹²	30-90	Therapist administers in standardised order	Paediatric health professionals early childhood specialists	Formal training not required. DVD, webinars and workshops available	Comprehensive manual/kit: £1089 Test kit provides most equipment
BOT-2 ¹³	40-60	Therapist administered in standardised order	Paediatric health professionals early childhood specialists	Formal training not required	Comprehensive manual/kit: £961 Test kit provides most equipment
MABC-2 ¹⁴	20-40	Therapist administers items in standardised order. Some flexibility allowed.	Research psychologists, OT, PT, Paediatricians	Formal training not required.	Comprehensive manual/ kit: £1191 Test kit provides most equipment
MAND ¹⁵	15-20	Therapist administers items in standardised order.	Professionals e.g. education, neurology, OT, PT, psychology etc.	Formal training not required.	Manual and test kit: £1366 includes equipment
NSMDA ¹⁶	20-45	Observation followed by therapist administration of test items.	PT, OT	Formal training not required (but is available)	Comprehensive manual: £35. Equipment not included
PDMS-2 ¹⁷	45-60 (20-30 for GM only)	Standardised procedure.	Paediatric health professionals, PE teachers, early intervention specialists	Formal training not required	Comprehensive manual/kit: £553 Includes some but not all equipment required
TGMD-2 ¹⁸	15-20	Standardised procedure.	Teachers, health professionals (OT, PT, doctors)	Formal training not required	Kit includes manual and record form: £128. Equipment not included

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition¹²; BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition¹³; MABC-2, Movement Assessment Battery for Children 2nd edition¹⁴; MAND, McCarron Assessment of Neuromuscular Development¹⁵; NSMDA, Neurological Sensory Motor Developmental Assessment¹⁶; PDMS-2, Peabody Developmental Motor Scales 2nd edition¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition¹⁸; GM, Gross motor; OT, Occupational Therapy; PT, Physiotherapy; PE, Physical Education

Table 3. Methodological quality of included articles

Test	First author, Year	Country	Population (Age, Diagnosis)	Internal consistency	Reliability	Measurement error	Content validity	Structural validity	Hypothesis testing	Cross- cultural validity	Criterion validity	Responsive - ness
BAYLEY III	Bayley ¹²	USA	1-42 mths	Fair	Fair	Good	Excellent	Good	Good	-	Good	-
	Spittle, et al. ⁴	Australia	2,4 yrs, Ex prem	-	-	-	-	-	-	-	Good	-
	Visser, et al. ²⁶	Netherlands	2.2-10.8 yrs, GDD, L.I.	-	-	-	Excellent	Poor	-	-	-	-
BOT-2	Wuang and Su ²⁷	Taiwan	4-12 yrs ID	Excellent	Excellent	Excellent	-	-	-	-	-	Fair
	Wuang, et al. ²⁸	Taiwan	3-6 yrs ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Bruininks and Bruininks ¹³	USA	4-21 yrs	Good	Fair (interrater) Fair (test-retest)	Good	Excellent	Good	-	-	Good	-
MABC-2 (AB 1)	Ellinoudis, et al. ²⁹	Greece	3-5.5 yrs	Excellent	Good	-	-	-	-	-	-	-
	Hua, et al. ²⁴	China	3-6 yrs	Excellent	Good	-	Excellent	Excellent	-	Poor	Excellent	-
	Logan, et al. ⁵	USA	3-6 yrs	-	-	-	-	-	Fair	-	Fair	-
	Smits-Engelsman, et al. ³⁰	Belgium	3-4 yrs	Poor	Poor	Poor	-	-	-	-	-	-
	Holm, et al. ³¹	Norway	7-9 yrs	-	Fair (interrater) Poor (intrarater)	Poor	-	-	-	-	-	-
MABC-2 (AB 2)	Kita, et al. ³²	Japan	7-10 yrs	Excellent	-	-	-	-	-	Poor	-	-
	Griffiths, et al. ³³	Australia	4-8 yrs	-	-	-	-	-	-	-	Good	-
MABC-2	Henderson, et al. ¹⁴	UK	3-16 yrs	-	Fair	Good	Excellent	-	-	-	-	-
	Niemeijer, et al. ²⁵	Netherlands + Belgium	-	-	-	-	-	-	-	Poor	-	-
	Schulz, et al. ³⁴	U.K	3-16 yrs	-	-	-	Excellent	Good	-	-	-	-
	Valentini, et al. ³⁵	Brazil	3-13 yrs	Fair	Fair	-	Fair	Poor	-	Poor	Poor	-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

	Wuang, et al. ²⁸	Taiwan	3-6 yrs, ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Wuang, et al. ³⁶	Taiwan	6-12 yrs DCD	Poor	Fair	Good	-	-	-	-	-	Fair
MAND	Hands, et al. ³⁷	Australia	10-17 yrs	-	-	-	-	Excellent	-	-	-	-
	McCarron ¹⁵	USA	7yrs	-	-	-	Fair	Poor	-	-	Poor	-
NSMDA	Danks, et al. ³⁸	Australia	2 + 4 yrs ELBW	-	-	-	-	-	-	-	Fair	-
	MacDonald and Burns ³⁹	Australia	2 + 4 yrs CP	-	-	-	-	Fair	-	-	Poor	-
	Burns, et al. ⁴⁰	Australia	1-24 mths VLBW	Poor	-	-	Poor	-	-	-	-	-
	Burns, et al. ⁴¹	Australia	1-mnths VLBW	-	-	-	-	Poor	-	-	Fair	-
PDMS-2	Hua, et al. ²⁴	China	3-6 yrs.	Excellent	Good	-	Excellent	Excellent	-	Poor	Excellent	-
	Wuang, et al. ²⁸	Taiwan	3-6 yrs ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Folio and Fewell ¹⁷	USA	0-71 mths	Good	-	Poor	Excellent	Good	Good	-	Poor	-
TGMD-2	Barnett, et al. ⁴²	Australia	4-8 yrs	-	Fair	-	-	-	-	-	-	-
	Farrokhi, et al. ⁴³	Iran	3-11 yrs	Fair	Fair	-	Fair	Fair	-	-	-	-
	Houwen, et al. ²¹	Netherlands	6-12 yrs VI	Fair	Fair	-	-	Fair	-	-	-	-
	Kim, et al. ⁴⁴	Korea	8-12 yrs ID	-	Poor	-	-	-	-	-	-	-
	Kim, et al. ⁴⁵	Korea	5-6 yrs	Poor	Fair	-	-	Poor	-	-	Poor	-
	Logan, et al. ⁵	USA	3-6 yrs	-	-	-	-	-	Fair	-	Fair	-
	Rudd, et al. ¹⁹	Australia	6-12 yrs	-	-	-	-	Good	-	-	-	-
	Simons, et al. ²³	Belgium	7-10 yrs ID	Good	Good (interrater) Poor (test-retest)	-	Excellent	Good	Good	-	-	-
	Valentini ²⁰	Brazil	3-10 yrs	Poor	Fair (test-retest) Good (intra,	-	Excellent	Good	-	Fair	Good	-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

interrater)											
Wong and Yin Cheung ²²	China	3-10 yrs	-	-	-	-	Fair	-	-	-	-
Ulrich ¹⁸	USA	3-10 yrs	Good	Fair (test-retest)	Fair	Poor	Good	-	-	Fair	-
Poor (interrater)											

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ Mths, Months; yrs, years; DCD, Developmental Coordination Disorder; VI, Vision Impairment; ID, Intellectual Disability; GDD, global developmental delay; L.I, Language Impairment; ELBW, Extremely Low Birth Weight; VLBW, Very Low Birth Weight; CP, Cerebral Palsy; prem, premature; USA, United States of America

For peer review only

1 Reliability

2 Internal consistency of assessments are summarised in Table 6. The BOT-2's high internal
3 consistency is well supported, including for children with an intellectual disability^{28 46}. The MABC-2
4 appears to have lower internal consistency than the BOT-2, which may relate to the limited number
5 of test items (eight) on the MABC-2. The highest values for internal consistency for the MABC-2
6 were obtained in specific populations (intellectual disability and developmental coordination
7 disorder) with poor to fair methodology only. Conversely the highest quality articles reported the
8 lowest values, although it should be noted that these assessed age band 1 (3-6 years) only. Internal
9 consistency is reported to be high for the PDMS-2, while the Bayley-III is shown to have excellent
10 internal consistency in children aged 24-42 months.. The TGMD-2 is reported by two good quality
11 (and four poor to fair quality) articles to have excellent internal consistency, including for children
12 with vision impairment and intellectual disability. The MAND is the only assessment tool included in
13 this review without published data of internal consistency or reliability in this age group.

14 The reliability findings are summarised in Table 6 and in Figures 2 and 3. Test-retest reliability was
15 excellent in the Bayley-III (Table 6), BOT-2 and PDMS-2; and was good to excellent in the MABC-2
16 and TGMD-2 (Figure 2). Intra-rater reliability was rarely investigated or reported for most tools, with
17 the TGMD-2 demonstrating better results than the MABC-2 (Figure 3). Only the TGMD-2 and
18 MABC-2 report inter-rater reliability values using an ICC (Figure 3)^{31 42}. Inter-rater reliability is also
19 supported in the BOT-2 with Pearson Correlation Coefficient and Kappa respectively. The studies
20 referred to in the test manuals for the TGMD-2, Bayley-III, BOT-2 and MABC-2 all report reliability
21 findings using Pearson's correlation, which is less ideal than an ICC or weighted kappa for statistical
22 analysis^{47 48}. Only studies reporting ICC's are visually represented in Figures 2 (test-retest) and 3
23 (inter and intra-rater). The TGMD-2 test-retest reliability results from Houwen, et al.²¹ were
24 believed to contain an error as the reported ICC was outside of the reported confidence intervals
25 (ICC 0.92, 0.82-0.91). This data set was therefore excluded from Figure 2.

26 Responsiveness was reported for the Bayley-III, BOT-2, MABC-2 and PDMS-2 with minimal
27 detectable change (MDC) or a standard error of measurement (SEM)²⁸. Sensitivity and specificity
28 for detecting change was shown to be satisfactory in the MABC-2, PDMS-2 and MABC-2²⁸ (Table 6).
29 There have been no studies to date on the responsiveness of the TGMD-2, NSMDA or MAND.

Table 4: Content and construct validity of assessment tools

Test	Content	Construct
BAYLEY III	Expert opinion for standard and low verbal version ^{12 26} . Literature reviews. Gross motor score correlated with Motor component 0.70 ¹²	Factor analysis. Difference in mean scores with pervasive developmental disorder, and specific language impairment ¹² . H_i (gross motor subset) = 0.52-0.97 for children with language impairment and 0.82-0.99 in control group ²⁶
BOT-2	Focus groups, product survey, pilot, national tryout and standardisation studies, professional reviews ¹³	Factor analysis, scores increase with age, discriminates between normal and children with DCD ($N=50$), high-functioning ASD ($N=45$) and mild-moderate ID ($N=66$) ¹³
MABC-2	Expert Panel, Stakeholder feedback, Literature review ³¹ Expert panel - clarity (validity content index 71.8-93.9, Kappa 0.76-0.88) and pertinence (98.5-99.3 and kappa 0.83-0.92) $p<0.001$ ³⁵	Factor analysis, correlation coefficients ²⁹ Subtest correlations 0.65-0.76 $p<0.001$. Discriminates between ASD and control group ³¹ . Structural equation modelling (for each age group) ³⁴ . Expert panel - adequate face validity ³⁵ . Significant difference between TD, DCD and at risk DCD scores ($\eta^2 = 0.63$) $p<0.0001$ ³⁵ . UK norms not appropriate to use with Dutch/Flemish children as under/over-estimate risk of motor impairment ²⁵ . In Chinese population: CFA initially rejected. Acceptable fit achieved after 2 items removed ²⁴ . Age band 2 shows good validity in Japanese population ³² .
MAND	Based on neuropsychological theory. Several rounds of revision/trials of tasks during development ¹⁵	Factor analysis ^{15 37} . Scores increase with age, and discriminate between typically developing children and those with head trauma or neurological dysfunction as well as gender ^{15 37}
NSMDA	Literature review. Developed by an experienced paediatric physiotherapist ⁴⁰	Factor analysis (up to 2 years of age) ^{40 41} . Stability of test results over time (up to 2 years) ^{40 41} .
PDMS-2	Literature review. Created by experts in the field. Revised with feedback from therapists guided revision. Hierarchical sequence of items ¹⁷	Item response modelling. Factor analysis. Differential item functioning analysis. Scores correlated with age ($r=0.80-0.93$) ¹⁷
TGMD-2	Expert Panel (3 PE teachers with post-grad qualifications) ¹⁸ . Translated version (Brazilian Portuguese) language clarity 0.96, pertinence >0.89 . Experts CVI for clarity and pertinence were also strong- $\alpha = 0.93$ clarity and $\alpha = 0.91$ pertinence ²⁰	Exploratory and confirmatory factor analysis ^{19 20 18 21 22 23} High and significant correlation of increasing age and increasing scores ⁴³ . Age and disability differentiation ^{18 23} Subtest correlation 0.41 ¹⁸ Gallop, running and leaping not well correlated with locomotion subscale. Object control significant & highly correlated ⁴⁵ . ANOVA - significant age effect for object control ²³ Moderate correlation between items and subset scores, and between subset scores and total score ²³

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸; H_i , scalability coefficient; CFA, Confirmatory Factor Analysis; TD, Typically Developing; ASD, Autism Spectrum Disorder, ID, Intellectual Disability; WPPSI, *Wechsler Preschool and Primary Scale of Intelligence*; WISC-R, Wechsler Preschool and Primary Scale of Intelligence-R; NDI, Neurodevelopmental Index; ANOVA, Analysis of Variance

Table 5: Criterion and predictive validity of assessment tools

Test	Criterion	Predictive
BAYLEY III	Given but mean age <22 months. Not relevant to study population. ¹²	Motor impairment at 4 years: Bayley III at 2 years <1SD = sensitivity 0.32-0.37 specificity 0.97 <2SD sensitivity 0.18-0.21 specificity 1.00. CP at 4 years: Bayley III at 2 years <1SD sensitivity 0.83 specificity 0.94. <2SD sensitivity 0.67 specificity 1.0 ⁴
BOT-2	MABC-2 $\rho = 0.92$ PDMS-2 $\rho = 0.88$ ($N = 38$) ²⁸ . PDMS-2 Total motor composite $r = 0.77$ ¹³ .	-
MABC-2	PDMS-2 $\rho = 0.631 - 0.84$ ^{28,24} . TGMD-2 $\rho = 0.45$ ⁵ . TGMD-2 standard scores ($r = 0.3, p < 0.02$) ³⁵ . BOT-2 $\rho = 0.90 - 0.92$ ²⁸ .	Classification groups (DCD, at risk and TD) remained same over time (6 months) $\chi^2 = 0.67 p = 0.72$ ³⁵ . Predictive of motor impairment over 6-12 months ($N=41$) ICC $0.88 p < 0.007$ ³⁵ . Scores at 4 years predictive of motor impairment at 8 years in children born <30 weeks gestation (PPV 79, sensitivity 79%, specificity 93%) ³³
MAND	Gross motor subscore: Low-moderate correlation with manual dexterity (-0.46 to 0.35), reaction time (-0.31 to -0.58), intelligence measures (WISC-R, Metropolitan Achievement Test) (0.30-0.39) and visual motor test (-0.33 to 0.39) ¹⁵	-
NSMDA	NSMDA at 2 years ($N = 148$) predictive of medical diagnosis $\chi^2 = 0.08 p = NS$ ⁴¹	Motor outcome at 11-13 yrs. NSMDA at 2 years - sensitivity 48.8%, specificity 82.4%, NSMDA at 4 years sensitivity 64.5%, and specificity 80%. PPV at 2 years 83% at 4 years 87% ³⁸ . If classified 'severe' at 24 months - approximately 50% chance walking at 4 years (moderate = 80%, mild = 93% minimal = 100%) ³⁹
PDMS-2	MABC-2 $\rho = 0.63 - 0.84$, ^{24,28} MABC-2 gross motor composite $\rho = 0.743$ ²⁴ BOT-2 $\rho = 0.88$ ²⁸ . Mullen Scales of Early Learning GMQ = 0.86 FMQ = 0.80 ¹⁷	-
TGMD-2	MABC-2 total $r = 0.49 p < 0.01$ ⁵ . 'Teacher report' $r = 0.34-0.45$. physical fitness $r = -0.47 - 0.55$ ⁴⁵ ($N=41$) Basic Motor Generalizations subtest of the CSSA $r = 0.63$. Locomotor 0.63 object control 0.41 ¹⁸	-

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ NS, Not Specified; SD, Standard Deviation; CP, Cerebral Palsy; TD, Typically Developing; ICC, Intraclass Correlation Coefficient; χ^2 , Chi Squared; NDI, Neurodevelopmental Index; CSSA, Comprehensive Scales of Student Abilities

Table 6: Reliability of assessment tools

Test	Internal Consistency	Test-Retest	Intra-rater	Inter-rater	Minimal detectable change	Minimal clinical important difference
BAYLEY III	GM $\alpha = 0.87$ -0.93 MC: $\alpha 0.90$ -0.96 (24-42 months) ¹²	Gross Motor subtest (N=47) $r=0.79$ Motor component $r=0.80$ ¹²	-	-	SEM Gross motor subtest 0.85-1.08. of Motor component = 3.00-4.74 (24-42 months) ¹²	-
BOT-2	(N = 100) $\alpha = 0.92$ ²⁷ (N = 141) $\alpha = 0.86$ ²⁸ 4-7 yrs (N= 620) $\alpha = 0.95$ 8-11 yrs (N= 450) $\alpha = 0.95$ ¹³	(N = 100) ICC = 0.99 ²⁷ (N = 141) ICC = 0.97 ²⁸ 4-7 yrs (N = 43) $r = 0.81$ (8-12 yrs (N= 44) $r = 0.80$ ¹³	-	Total motor composite 4-21 yrs (N = 47) $r = 0.98$ ¹³	4.18 (sensitivity 55.10% specificity 72.55%) ²⁷ 7.43 (sensitivity 42.49% specificity 65.72%) ²⁸	6.53 (sensitivity 48.98% specificity 76.47%) ²⁷ 6.55 (sensitivity 49.99% specificity 58.78%) ²⁸
MABC-2 (AB 1)	(N = 60) M.D $\alpha = 0.51$, A&C $\alpha = 0.70$, Bal $\alpha = 0.66$ ²⁹ (N = 1823) $\alpha = 0.502$ ²⁴ (N=50) $\alpha = 0.81$ -0.87 ³⁰	(N=60) ICC = 0.85 ²⁹ Item ICC's 0.830-0.985 ²⁴ ICC test-retest = 0.83 ³⁰ Inter-rater test-retest ICC = 0.79 ³⁰	(N=28) $\kappa = 0.71$ ³⁰	Item ICC's range 0.892-0.998 ²⁴ (N=22) $\kappa = 0.60$ ³⁰	(N=28) Intrarater MDC = 3.43 (N=22) Inter-tester MDC = 3.81 ³⁰	-
MABC-2 (AB 2)	Translated version (Japanese) (N=132) $\alpha = 0.602$ ³²	-	ICC = 0.64 ³¹	ICC 0.63 ³¹	Intra-rater SDC TTS: +/- 11.7 TSS +/- 3.3. Inter-rater SDC TTS +/-16.0 TSS +/- 3.8 ³¹	-
MABC-2	Subscales $\alpha = 0.78$ (M.D = 0.77, BS = 0.52, Bal = 0.77) ³⁵ $\alpha = 0.88$ ³⁶ (N = 141) $\alpha = 0.88$ ²⁸	N=60 (all 3 age bands) $r=0.80$ ¹⁴ $r=0.74$ $p<0.0001$ (standard score). ICC standard score = 0.85 ³⁵ ICC 0.96 ³⁶ N = 141 ICC =0.96 ²⁸	ICC 0.88 ³⁵	ICC 0.96-0.99 ³⁵	SEM 1.34 (95%CI) = 3 ¹⁴ 1.83 (95%CI) ³⁶ 1.83 (sensitivity 69.69% specificity 52.10%) ²⁸	1.39 (sensitivity 72.47% specificity 46.18%) ^{28 36}
MAND	-	-	-	-	-	-
NSMDA	Cross correlation matrix Item scoring (12+24months) 0.73 $p<0.001$, Functional grade (12+24months) 0.87 $p<0.001$ ⁴⁰	-	-	-	-	-
PDMS-2	(N=141) $\alpha=0.89$ ²⁸ 24-35m $\alpha=0.97$, 36-47m $\alpha=0.95$, 48-59m $\alpha=0.97$, 60-71m $\alpha=$	N=141 ICC= 0.97 ²⁸	unable to extract data for ≥ 24 months	unable to extract data for ≥ 24 months ¹⁷	7.76 (sensitivity 60.65% specificity 74.13%) ²⁸ SEM 24-	8.39 (sensitivity 61.65% specificity 71.34%) ²⁸

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

	o.98. For subgroups† $\alpha=0.99$ ¹⁷		¹⁷	59 months = 3, 60-71m = 2 ¹⁷
TGMD-2	($N=1438$) $\alpha=0.80$ ⁴³ $N=75$ Locomotor subset $\alpha=0.71$ object control $\alpha=0.72$ ²¹ $N=120$ $\alpha = 0.72$ ⁴⁵ $N=99$ $\alpha = 0.90$ ²³ $N = 1208$ Cronbach's $\alpha = 0.91$ (gross motor quotient). Locomotor 0.85 and object control 0.88. Note SEM GMQ = 4-5 SEM subsets=1 ¹⁸	$N=63$ ICC=0.81 95% CI ⁴³ $N=23$ ICC=0.92 total 95% CI ²¹ $N=99$ $r=0.98$ ²³ Locomotor test $r = 0.90$ $p < 0.0001$ object control test $r = 0.91$ $p < 0.001$ ²⁰ $N = 75$ $r=0.96$ overall (3-5 yrs $r = 0.91$), 6-8 years $r = 0.95$), (9-10 years $r = 0.94$) ¹⁸	$N=32$ ICC=0.97 95% CI ⁴³ $N=25$ ICC=0.95 95% CI ²¹ ICC = 0.78 ⁴⁴ ICC=0.92-0.99 ²⁰	Obj ICC=0.93 ⁴² ($N=50$) ICC=0.89 ²¹ ICC=0.75 ⁴⁴ $N=8$ $r= 1.00$ ²³ L.S ICC=0.88 Obj ICC=0.89 ²⁰ $N = 30$ $r=0.98$ ¹⁸

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ GM, Gross Motor Subset; MC, Motor Component; K, Kappa Coefficient; M.D, Manual Dexterity; BS, Ball Skills; BAL, Balance; A&C, Aiming and catching; SDC, Smallest Detectable Change; TTS, Total Test Score; TSS, Total Standard Score; †, gender, ethnicity, speech/language or physical disorder; Obj, Object Control Subset; L.S, Locomotion Subset

Discussion

This review identified seven gross motor assessment tools appropriate for use in clinical or research settings, each with their own strengths and limitations. Interestingly, only one of the seven assessments (TGMD-2) measured gross motor skills in isolation. This is likely a reflection on current practice to assess children's development as a whole, rather than assessing individual domains in isolation. A gross motor assessment embedded within a developmental assessment, such as that of the Bayley-III may be more appropriate than an isolated gross motor assessment for children where there is suspicion of multiple impairments.

A review by Slater, et al.⁸ reported that the TGMD-2 and the MABC (first edition) were recommended for assessing gross motor skills in children with developmental coordination disorder, but found that the MABC needed further evidence of validity. Cools, et al.⁴⁹ also published a detailed review of the clinical utility of gross motor assessment tools for children, but did not address the validity, reliability or responsiveness to change of these measures. This review adds to the literature by including updated information on the psychometric properties of the measures and a thorough methodological assessment using the COSMIN checklist which allows the reader to interpret these results with confidence. We have identified ten additional publications to support the content, construct and criterion validity of the MABC-2 and have demonstrated an overall higher methodological quality of the papers assessing the MABC-2 when compared with the TGMD-2. Papers that had been given lower methodological scores on the COSMIN can be attributed to inadequate reporting statistical methods, small sample sizes and non-independent assessors. Further research in this area should consider addressing these limitations in their study design to reduce potential error and increase confidence when interpreting results.

Content validity has been established for five of the included assessment tools, however, further research into the content validity for the MAND and NSMDA is required. The NSMDA's ability to predict a diagnosis of CP and motor outcomes over time does support its content validity, however the methodology scored as poor to fair on the COSMIN and as such content validity cannot be fully established. The use of expert panels, focus groups and/or stakeholder feedback for the BOT-2, MABC-2, TGMD-2 and PDMS-2 demonstrate thorough consideration of the relevance and comprehensiveness of the each test's assessment items during development.

1 The TGMD-2 is the only assessment tool considered to have well established construct validity, with
2 several papers reporting factor analysis. The NSMDA has undergone factor analysis for children up to, but
3 not beyond two years of age and as such further research is needed to support its validity in older children.
4 All other included assessment tools have undergone factor analysis assessment of their construct validity in
5 one paper and are supported by the ability to discriminate between medical diagnosis or age, and as such
6 are considered to have adequate construct validity. The criterion validity indicates that the TGMD-2 may
7 be measuring a slightly different construct to the other assessment tools included in this study as it has
8 poor agreement with the MABC-2, which in turn has good agreement with the PDMS-2 and the BOT-2.
9 This difference may be related to the inclusion of the assessment of quality of movement in the TGMD-2,
10 or the inclusion of balance and/or fine motor tasks on the other assessments. There is scope to investigate
11 the criterion validity of the MAND and the gross motor subsections of the Bayley-III and the NSMDA with
12 the other assessment tools in this study in the future.

13 The BOT-2 was the only assessment tool to have its reliability assessed with excellent methodology. In
14 conjunction with its reported results it can be considered to have the strongest evidence for internal
15 consistency and test-retest reliability out of the included assessment tools. The PDMS-2 and the MABC-2
16 can be considered to have the next best established test-retest reliability with good methodological
17 quality. The reported test-retest reliability values for the TGMD-2 are impacted by the poor to fair
18 methodological quality, and further high quality research needs to be done to support its body of evidence.
19 Test-rest, inter or intra-rater reliability has not been assessed in the MAND and NSMDA. In the clinical
20 context gross motor assessments are often repeated over time or between therapists and as such these
21 measures of reliability should be established. The Bayley-III would also benefit from further research into
22 its reliability, with no published inter or intra-rater reliability measures, and with only one, fair quality
23 report of good test-retest reliability.

24 As yet there is little evidence to support the use of these assessments as outcome measures. The inclusion
25 in some of the articles of minimal detectable change (MDC) and minimal clinically important difference
26 (MCID) is valuable for clinicians⁷. The difference between MDC and MCID is also of importance, as a change
27 in score does not necessarily relate to a meaningful change for the child or their family. Only the Bayley,
28 BOT-2, MABC-2 and PDMS-2 have a reported MCID with satisfactory sensitivity and specificity, however,
29 due to the fair methodological quality used to obtain these values they cannot be utilised with a high level
30 of confidence until further studies have been performed. The TGMD-2 was created in part to be used as an
31 outcome measure, however there are no articles to date investigating its responsiveness to change¹⁸. It
32 should also be noted that all of the included assessment tools measure impairment and activity limitations,
33 but do not specifically address the other elements of the International Classification of Functioning,

1 Disability and Health (ICF) domains of participation, personal factors and environment². Clinicians should
2 utilise appropriate assessments or questionnaires to ensure that these domains of health are also
3 addressed in line with World Health Organisation guidelines².

4 When considering a test's reliability all three elements of test error should be taken into account – these
5 can be described as time sampling (assessed with test-retest reliability), content sampling (assessed as
6 internal consistency), and inter-scorer difference (or interrater reliability)¹⁸. This is one of the reasons that
7 clinicians should consider repeating assessments and/or completing a second alternative assessment. All
8 assessments should be interpreted in conjunction with clinical reasoning and observation. Included
9 assessment tools are not intended to be diagnostic on their own; results need to be combined with other
10 assessments and expert opinion to arrive at a clinical diagnosis.

11 All of the included assessment tools were found to have merits and limitations in their clinical utility the
12 body of evidence to support their use. Clinicians and researches should select their assessment tool with
13 consideration of psychometric properties (inclusive of the methodological rigour behind them), clinical
14 utility and for the population, situation and age group in question.

15 A potential limitation of this study was that one author screened the titles and abstracts, which may have
16 led to a sampling bias. Whilst care was taken to include all potentially relevant papers and assessment
17 tools until the second round of assessment with two authors, the potential for exclusion of papers relevant
18 to this review remains. A second limitation was the restriction of included papers and manuals to those
19 published in English. Unfortunately this resulted in the exclusion of three assessment tools that have been
20 reported as commonly used in Europe: The Motoriktest für Vier- bis Sechjährige Kinder (MOT 4-6), the
21 Körperkoordinationstest für Kinder (KTK) and the Maastrichtse Motoriek Test (MMT)⁴⁹. The authors also
22 note the third edition of the TGMD is soon to be published and will need to be subjected to a similar level
23 of assessment of psychometric properties in the future.

24 Clinicians and parents who need guidance to set realistic therapy goals and to understand future
25 intervention requirements benefit from understanding a test's predictive ability. The NSMDA and the
26 MABC-2 are the only tools that have demonstrated long term (≥4 years follow up) predictive validity, while
27 the Bayley-III has good predictive validity at 2 years for future movement difficulties and for the diagnosis
28 of cerebral palsy at 4 years. However, further research into the long-term predictive validity of all included
29 gross motor assessment tools is warranted.

30 While validity and reliability should guide selection of assessment tools, clinical utility must also be taken
31 into consideration. Most tests have ongoing costs associated with forms and equipment replacement,

1 which may be prohibitive to some users. The NSMDA requires the therapist to handle the child for several
2 items which should be considered in relation to manual handling policies of institutions. Assessment
3 burden for children and families should also be taken into consideration when selecting an assessment
4 tools. Younger children are more likely to be distracted and may not understand test items as well, which
5 may also increase assessment times³⁰.

6 When a new edition of an assessment tools is released resulting in a change in age groups, scoring or tasks
7 it is insufficient to rely on the psychometric assessments that were performed on the original test. The
8 MABC-2 manual provides justification for the inclusion of reliability and validity assessment of the original
9 MABC¹⁴, however, owing to the significant changes in age groups and tasks between editions these were
10 not included for the analysis of the MABC-2 in this review. Two studies quoted in the MABC-2 manual to
11 support the validity and reliability are both unpublished works and as such are also unable to be included
12 in this systematic review. This could indicate a publication for the MABC-2.

13 The thorough methodological assessment of the included articles using the COSMIN checklist should be
14 seen as a strength of this paper, as should the range of assessment tools included in this review. While it
15 has previously been argued that the 'worst score counts' criteria in the COSMIN creates a floor effect⁵⁰,
16 the COSMIN authors argue that only 'fatal flaws' contribute to an overall score of poor¹⁰. There are few
17 tools available to assess the psychometric properties of assessment tools and arguably none so robustly
18 validated as the COSMIN.

19 There are many appropriate gross motor assessment tools available for use in research and clinical settings
20 today. Most of the available tools demonstrate adequate validity and reliability in children aged 2-12 and
21 as such the authors do not believe that new assessment tools need to be developed for use. There is scope
22 however to improve the evidence of inter and intra-rater reliability and predictive validity should be
23 ascertained over a longer period of time and with greater methodological rigour. Tools also need clearer
24 assessment of their responsiveness to change to assist clinicians and researchers with outcome measure
25 selection. Researchers should be mindful of the methods they use to assess validity and reliability. Clarity
26 of reporting, statistical methods and sample sizes should be carefully considered to ensure the highest
27 quality of evidence.

28 Conclusion

29 Currently available gross motor assessment tools for children have good to excellent content and construct
30 validity. The BOT-2, MABC-2, PDMS-2 and TGMD-2 are the most reliable assessments in this age group. The

1 Bayley-III has the best predictive validity at 2 years of age, and the NSMDA and the MABC-2 both have
2 good predictive validity at 4 years of age. There is scope for further research into the predictive validity,
3 reliability and responsiveness of gross motor assessment tools in preschool and school aged children. In
4 practice clinicians should choose assessments with consideration of their psychometric properties in the
5 context of the child that they are assessing.
6
7
8
9
10
11

12 **Author Contributions**

13
14 8 All individuals listed as authors meet the appropriate authorship criteria and have approved the
15 acknowledgement of their contributions. AG was responsible for the drafting of the paper and liaising with
16 9 the co-authors on findings and conclusions. RT contributed to the paper through interpretation of data,
17 10 completing methodological assessments and revising manuscript content throughout its development.
18 11 A/Profs PEM and AJS both contributed to the paper through assisting with the development of research
19 12 design, interpretation of data and revising manuscript content through its development.
20
21
22
23
24
25
26
27

28 **Data Sharing Statement**

29
30 16 This paper includes data obtained from reviewing papers of published manuscripts. Data can be accessed
31 by contacting the primary author.
32
33
34
35

36 **Figures**

37
38
39 19 Figure 1. PRISMA flow diagram detailing study selection
40

41 20 Figure 2. Test re-test reliability of gross motor assessment tools
42
43

44 21 Figure 2 legend: BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition ¹³; MABC-2, Movement
45 22 Assessment Battery for Children 2nd edition ¹⁴; PDMS-2, Peabody Developmental Motor Scales 2nd edition
46 23 ¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition ¹⁸.
47
48
49

50 24 Figure 3. Inter and interrater reliability of gross motor assessment tools
51

52 25 Figure 3 legend: MABC-2, Movement Assessment Battery for Children 2nd edition ¹⁴; TGMD-II, Test of Gross
53 26 Motor Development 2nd edition ¹⁸
54
55
56
57
58
59

1 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

References

1. Piek JP, Baynam GB, Barrett NC. The relationship between fine and gross motor ability, self-perceptions and self-worth in children and adolescents. *Hum Mov Sci* 2006;25(1):65-75. doi: <http://dx.doi.org/10.1016/j.humov.2005.10.011>
2. World Health Organization. International Classification of Functioning, Disability and Health: ICF: World Health Organization 2001.
3. Magalhaes LC, Cardoso AA, Missiuna C. Activities and participation in children with developmental coordination disorder: a systematic review. *Res Dev Disabil* 2011;32(4):1309-16. doi: 10.1016/j.ridd.2011.01.029 [published Online First: 2011/02/19]
4. Spittle AJ, Spencer-Smith MM, Eeles AL, et al. Does the Bayley-III Motor Scale at 2 years predict motor outcome at 4 years in very preterm children? *Developmental Medicine And Child Neurology* 2013;55(5):448-52. doi: 10.1111/dmcn.12049
5. Logan SW, Robinson LE, Getchell N. The Comparison of Performances of Preschool Children on Two Motor Assessments. *Perceptual and Motor Skills* 2011;113(3):715-23.
6. Venetsanou F, Kambas A. Environmental factors affecting preschoolers' motor development. *Early Childhood Education Journal* 2010;37(4):319-27.
7. Mokkink LB, Terwee CB, Patrick DL, et al. The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. *J Clin Epidemiol* 2010;63(7):737-45. doi: 10.1016/j.jclinepi.2010.02.006 [published Online First: 2010/05/25]
8. Slater LM, Hillier SL, Civetta LR. The clinimetric properties of performance-based gross motor tests used for children with developmental coordination disorder: A systematic review. *PEDIATRIC PHYSICAL THERAPY* 2010;22(2):170.
9. Brown T, Lalor A. The Movement Assessment Battery for Children—Second Edition (MABC-2): A Review and Critique. *Phys Occup Ther Pediatr* 2009;29(1):86-103.
10. Terwee CB, Mokkink LB, Knol DL, et al. Rating the methodological quality in systematic reviews of studies on measurement properties: a scoring system for the COSMIN checklist. *Qual Life Res* 2012;21(4):651-7. doi: 10.1007/s11136-011-9960-1 [published Online First: 2011/07/07]
11. Law M. Outcome measures rating form. Ontario, Canada: CanChild Centre for Disability Research, 2004.
12. Bayley N. Bayley Scales of Infant Development and Toddler Development: Technical Manual: The PsychCorp 2006.
13. Bruininks R, Bruininks B. Bruininks-Oseretsky Test of Motor Proficiency—2nd Edition (BOT-2): Manual. Circle Pines: MN: AGS Publishing 2005.
14. Henderson SE, Sugden DA, Barnett AL. Movement assessment battery for children-2: Movement ABC-2: Examiner's manual: Pearson 2007.
15. McCarron LT. MAND: McCarron assessment of neuromuscular development, fine and gross motor abilities: McCarron-Dial Systems, Incorporated 1997.

- 1 16. Burns YR. N.S.M.D.A Physiotherapy Assessment for Infants and Young Children Second Edition. Brisbane,
2 Queensland: CopyRight Publishing Company 2014.
- 3 3 17. Folio M, Fewell R. Peabody Developmental Motor Scales. Examiner's Manual . 2nd Edition. Austin, Texas.:
4 Pro-Ed. 2000.
- 5 4
6 5 18. Ulrich DA. Test of gross motor development-2. *Austin: Prod-Ed 2000*
- 7 5
8 6 19. Rudd J, Butson ML, Barnett L, et al. A holistic measurement model of movement competency in children.
9 *Journal of Sports Sciences 2016;34(5):477-85.*
- 10 7
11 8 20. Valentini NC. Validity and reliability of the TGMD-2 for Brazilian children. *Journal of Motor Behavior*
12 2012;44(4):275-80.
- 13 9
14 10 21. Houwen S, Hartman E, Jonker L, et al. Reliability and validity of the TGMD-2 in primary-school-age children
15 with visual impairments. *Adapted Physical Activity Quarterly 2010;27(2):143-59.*
- 16 11
17 12 22. Wong KYA, Yin Cheung S. Confirmatory factor analysis of the Test of Gross Motor Development-2.
18 *Measurement in Physical Education & Exercise Science 2010;14(3):202-09. doi:*
19 13 *10.1080/10913671003726968*
- 20 14
21 15 23. Simons J, Daly D, Theodorou F, et al. Validity and reliability of the TGMD-2 in 7-10-year-old Flemish
22 children with intellectual disability. *Adapted physical activity quarterly : APAQ 2008;25(1):71-82.*
- 23 16
24 17 24. Hua J, Gu G, Meng W, et al. Age band 1 of the Movement Assessment Battery for Children-Second Edition:
25 exploring its usefulness in mainland China. *Research in Developmental Disabilities 2013;34(2):801-8.*
- 26 18
27 19 25. Niemeijer AS, van Waelvelde H, Smits-Engelsman BC. Crossing the North Sea seems to make DCD
28 disappear: cross-validation of Movement Assessment Battery for Children-2 norms. *Hum Mov Sci*
29 2015;39:177-88. doi: 10.1016/j.humov.2014.11.004
- 30 20
31 21 26. Visser L, Ruiters SAJ, Van der Meulen BF, et al. Low verbal assessment with the Bayley-III. *Research in*
32 22 *Developmental Disabilities 2015;36:230-43.*
- 33 23
34 24 27. Wuang YP, Su CY. Reliability and responsiveness of the Bruininks-Oseretsky Test of Motor Proficiency-
35 Second Edition in children with intellectual disability. *Research in Developmental Disabilities 2009;30(5):847-*
36 25 *55.*
- 37 26
38 27 28. Wuang YP, Su CY, Huang MH. Psychometric comparisons of three measures for assessing motor functions
39 in preschoolers with intellectual disabilities. *Journal of Intellectual Disability Research 2012;56(6):567-78.*
- 40 28
41 29 29. Ellinoudis T, Evaggelina C, Kourtessis T, et al. Reliability and validity of age band 1 of the Movement
42 Assessment Battery for Children--second edition. *Res Dev Disabil 2011;32(3):1046-51. doi:*
43 30 *<http://dx.doi.org/10.1016/j.ridd.2011.01.035>*
- 44 31
45 32 30. Smits-Engelsman BCM, Niemeijer AS, van Waelvelde H. Is the Movement Assessment Battery for Children-
46 2nd edition a reliable instrument to measure motor performance in 3 year old children? *Research in*
47 33 *Developmental Disabilities 2011;32(4):1370-77.*
- 48 34
49 35 31. Holm I, Tveter AT, Aulie VS, et al. High intra- and inter-rater chance variation of the movement assessment
50 battery for children 2, ageband 2. *Research in Developmental Disabilities 2013;34(2):795-800.*
- 51 36
52
53
54
55
56
57
58
59
60

- 1 32. Kita Y, Suzuki K, Hirata S, et al. Applicability of the Movement Assessment Battery for Children-Second
2 Edition to Japanese children: A study of the Age Band 2. *Brain & Development* 2016;38(8):706-13.
- 3 33. Griffiths A, Morgan P, Anderson PJ, et al. Predictive value of the Movement Assessment Battery for
4 Children - Second Edition at 4 years, for motor impairment at 8 years in children born preterm. *Dev Med*
5 *Child Neurol* 2017;59(5):490-96. doi: 10.1111/dmcn.13367 [published Online First: 2017/01/10]
- 6 34. Schulz J, Henderson SE, Sugden DA, et al. Structural validity of the Movement ABC-2 test: factor structure
7 comparisons across three age groups. *Research in Developmental Disabilities* 2011;32(4):1361-9.
- 8 35. Valentini NC, Ramalho MH, Oliveira MA. Movement assessment battery for children-2: translation,
9 reliability, and validity for Brazilian children. *Res Dev Disabil* 2014;35(3):733-40. doi:
10 <http://dx.doi.org/10.1016/j.ridd.2013.10.028>
- 11 36. Wang YP, Su JH, Su CY. Reliability and responsiveness of the Movement Assessment Battery for Children-
12 Second Edition Test in children with developmental coordination disorder. *Developmental Medicine & Child*
13 *Neurology* 2012;54(2):160-5.
- 14 37. Hands B, Larkin D, Rose E. The psychometric properties of the McCarron Assessment of Neuromuscular
15 Development as a longitudinal measure with Australian youth. *Human Movement Science* 2013;32(3):485-
16 97.
- 17 38. Danks M, Maideen MF, Burns YR, et al. The long-term predictive validity of early motor development in
18 "apparently normal" ELBW survivors. *Early Human Development* 2012;88(8):637-41.
- 19 39. MacDonald J, Burns Y. Performance on the NSMDA During the First and Second Year of Life to Predict
20 Functional Ability at the Age Of 4 in Children with Cerebral Palsy. *Hong Kong Physiotherapy Journal*
21 2005;23(1):40-45. doi: 10.1016/S1013-7025(09)70058-2
- 22 40. Burns YR, Ensbey RM, Norrie MA. The Neuro-sensory motor developmental assessment part 1:
23 development and administration of the test. *Aust J Physiother* 1989;35(3):141-49.
- 24 41. Burns YR, Ensbey RM, Norrie MA. The neuro-sensory motor developmental assessment part II: predictive
25 and concurrent validity. *Aust J Physiother* 1989;35(3):151-57.
- 26 42. Barnett LM, Minto C, Lander N, et al. Interrater reliability assessment using the Test of Gross Motor
27 Development-2. *Journal of Science and Medicine in Sport* 2014;17(6):667-70. doi:
28 <http://dx.doi.org/10.1016/j.jsams.2013.09.013>
- 29 43. Farrokhi A, Zareh Zadeh M, Karimi Alvar L, et al. Reliability and validity of test of gross motor development-
30 2 (Ulrich, 2000) among 3-10 aged children of Tehran City. *Journal of Physical Education and Sports*
31 *Management* 2014;5(2):18-28. doi: 10.5897/JPEM12.003
- 32 44. Kim Y, Park I, Kang M. Examining rater effects of the TGMD-2 on children with intellectual disability.
33 *Adapted Physical Activity Quarterly* 2012;29(4):346-65.
- 34 45. Kim CI, Han DW, Park IH. Reliability and validity of the test of gross motor development-II in Korean
35 preschool children: applying AHP. *Research in Developmental Disabilities* 2014;35(4):800-7.
- 36 46. Wang YP, Lin YH, Su CY. Rasch analysis of the Bruininks-Oseretsky Test of Motor Proficiency-Second
37 Edition in intellectual disabilities. *Research in Developmental Disabilities* 2009;30(6):1132-44.

- 1 47. Spittle AJ, Doyle LW, Boyd RN. A systematic review of the clinimetric properties of neuromotor
2 assessments for preterm infants during the first year of life. *Dev Med Child Neurol* 2008;50(4):254-66. doi:
3 10.1111/j.1469-8749.2008.02025.x [published Online First: 2008/01/15]
4
5 48. McDowell I. Measuring health: a guide to rating scales and questionnaires: Oxford university press 2006.
6
7 49. Cools W, De Martelaer K, Samaey C, et al. Movement skill assessment of typically developing preschool
8 children: a review of seven movement skill assessment tools.(Report). *Journal of Sports Science and Medicine*
9 2009;8(2):154.
10
11 50. Adair B, Said CM, Rodda J, et al. Psychometric properties of functional mobility tools in hereditary spastic
12 paraplegia and other childhood neurological conditions. *Dev Med Child Neurol* 2012;54(7):596-605.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

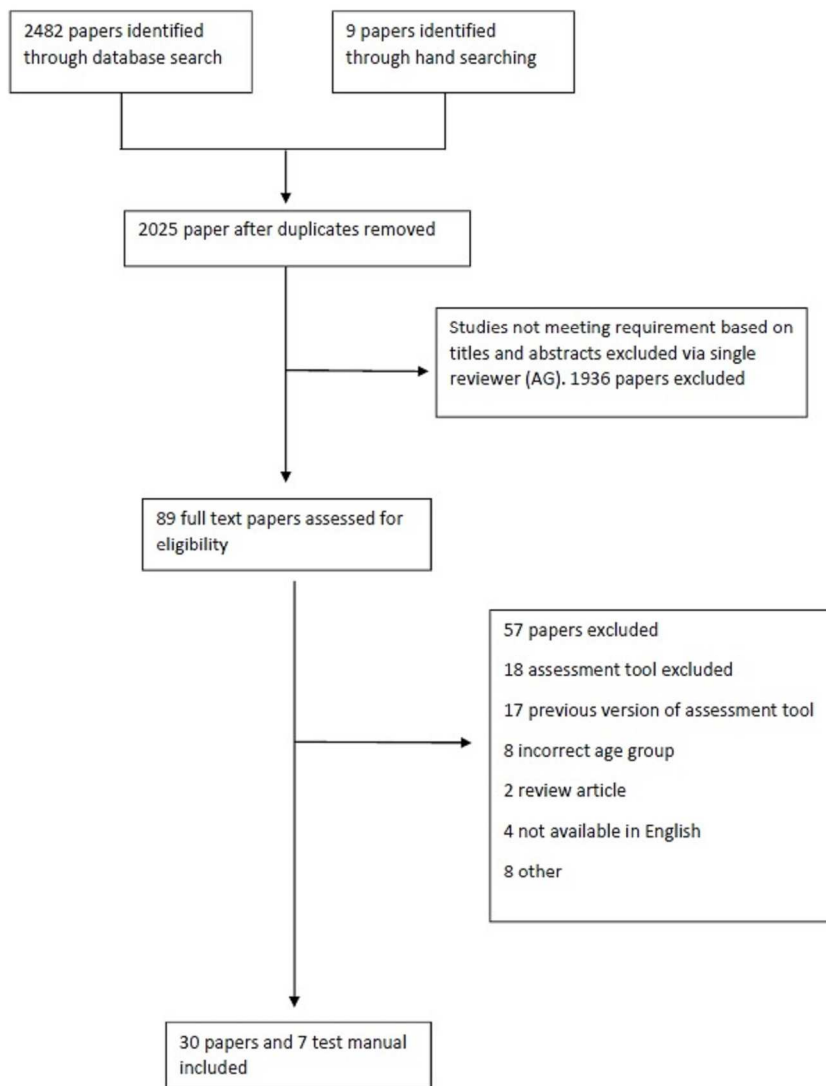


Figure 1. PRISMA flow diagram detailing study selection

195x256mm (300 x 300 DPI)

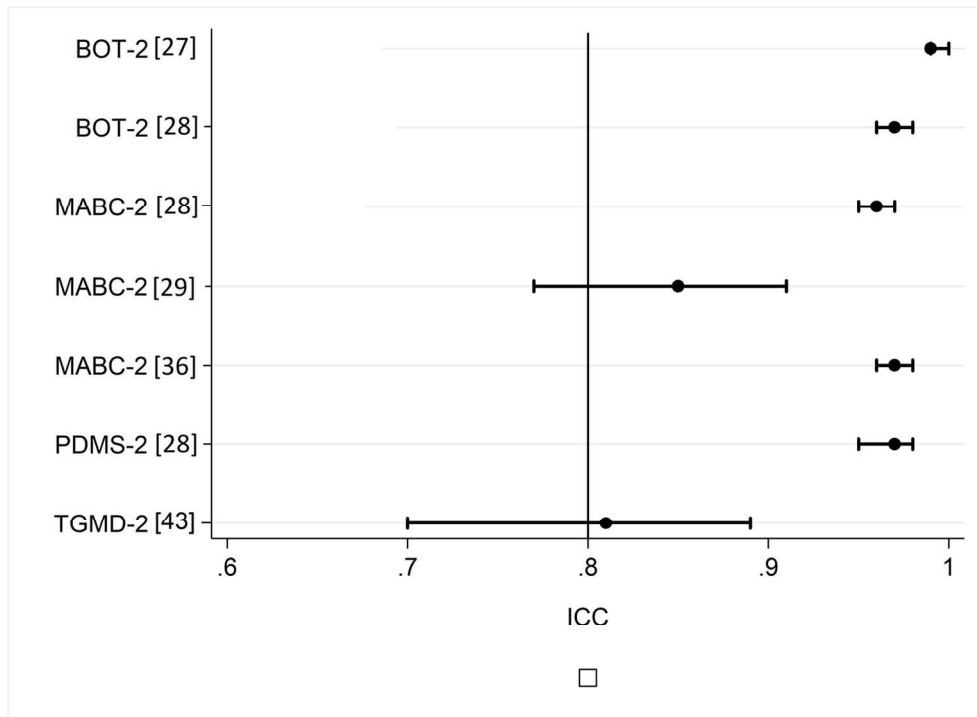


Figure 2. Test re-test reliability of gross motor assessment tools
Figure 2 legend: BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition 13; MABC-2, Movement Assessment Battery for Children 2nd edition 14; PDMS-2, Peabody Developmental Motor Scales 2nd edition 17; TGMD-II, Test of Gross Motor Development 2nd edition 18.

139x102mm (300 x 300 DPI)

Only

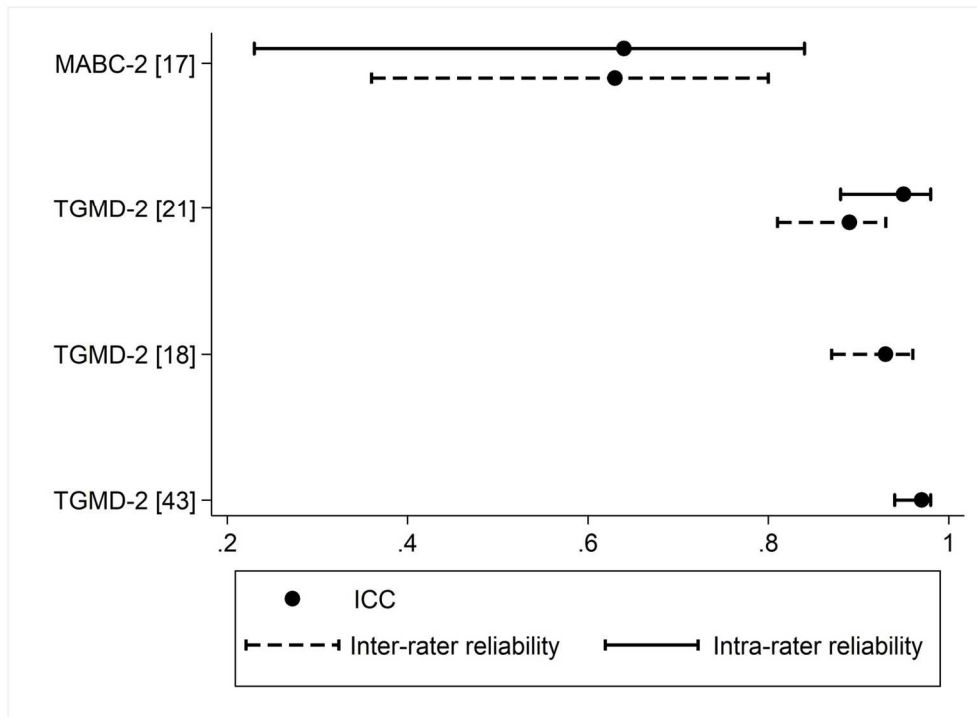


Figure 3. Inter and interrater reliability of gross motor assessment tools
 Figure 3 legend: MABC-2, Movement Assessment Battery for Children 2nd edition 14; TGMD-II, Test of Gross Motor Development 2nd edition 18

139x102mm (300 x 300 DPI)

Supplementary Table 1. OVID Medline database search (1996 to present)

Search No	Search	Yield 5/5/17
1	Child/	811722
2	Child, Preschool/	457484
3	paediatric*.mp.	45528
4	Motor Skills/	12726
5	Motor Activity/	64838
6	gross motor.mp.	3821
7	Psychomotor Disorders/	2609
8	Motor Skills Disorders/	2580
9	Developmental Disabilities/	13484
10	developmental coordination disorder.mp.	845
11	Movement/ph (physiology)	22342
12	Questionnaires/	336296
13	"Outcome Assessment (Health Care)"/	57491
14	scale*.mp.	608566
15	instrument*.mp.	197131
16	outcome*.mp.	1813266
17	measure*.mp.	2255187
18	evaluat*.mp.	2552240
19	assess*.mp.	2273012
20	"Task Performance and Analysis"/	22017 (or 5969)
21	Reproducibility of Results"/	319899
22	1 or 2 or 3	936097
23	4 or 5 or 6 or 7 or 8 or 9 or 10 or 11	116200
24	12 or 13 or 14 or 15 or 16 or 17 or 18 or 19 or 20	6579985
25	21 and 22 and 23 and 24	1152

Supplementary table 2: CINAHL plus database search

Search Number	Search	Yield
		2/7/17
S1	(MH "Child")	338732
S2	(MH "Child, Preschool")	156847
S3	"paediatric"	14351
S4	(MH "Motor Skills")	7420
S5	(MH "Motor Activity")	9664
S6	(MH "Psychomotor Performance")	9457
S7	(MH "Motor Skills Disorders")	1515
S8	(MH "Developmental Disabilities")	7114
S9	(MH "Child Development Disorders")	1708
S10	"gross motor"	2234
S11	(MH "Clinical Assessment Tools")	110291
S12	(MH "Outcome Assessment")	29335
S13	(MH Physical Therapy Assessment")	2027
S14	"scale"	290924
S15	instrument*	113984
S16	outcome*	602132
S17	measure*	534465
S18	evaluat*	760447
S19	assess*	710533
S20	(MH Reliability and Validity")	11043
S21	S1 OR S2 OR S3	387809
S22	S4 OR S5 OR S6 OR S7 OR S8 OR S9 OR S10	35562
S23	S11 OR S12 OR S13 OR S14 OR S15 OR S16 OR S17 OR S18 OR S19	1763829
S24	S20 AND S21 AND S22 AND S23	144

Supplementary table 3: EMBASE database search (1974-present)

Search No	Search	Search Yield 5/5/17
1	Child/	1346641
2	Preschool child/	457981
3	paediatric*.mp.	86036
4	Motor performance/	57571
5	Motor Activity/	39751
6	psychomotor performance	19515
7	Motor development/	4906
8	Motor dysfunction/	53155
9	Developmental disorder/	30473
10	Gross motor.mp.	6840
11	Outcome Assessment/	358121
12	Outcome measure.mp.	60507
13	Questionnaire/	513199
14	Task performance/	125167
15	Functional assessment/	55415
16	Clinical assessment tool/	19865
17	evaluat*.mp.	3874341
18	instrument*.mp.	515930
19	outcome*.mp.	2425627
20	Assess*.mp.	3815907
21	Scale*.mp.	903216
22	Measure*.mp.	3444366
23	Measurement accuracy/	18209
24	Measurement repeatability/	2849
25	Reproducibility/	173988
26	Validity/	40192
27	Reliability/	114002
28	1 or 2 or 3	1535605
29	4 or 5 or 6 or 7 or 8 or 9 or 10	195237
30	11 or 12 or 13 or 14 or 15 or 16 or 17 or 18 or 19 or 20 or 21 or 22	10121631
31	23 or 24 or 25 or 26 or 27	324779
32	28 and 29 and 30 and 31	1105

Supplementary table 4. Allied and Complementary Medicine Database (AMED) database search: (1985-present)

Search Number	Search	Yield 2/7/17
1	Child/	15192
2	Child preschool/	1223
3	Adolescent/	3979
4	paediatric*.mp.	812
5	1 or 2 or 3 or 4	18429
6	Motor skills/	1220
7	Motor activity/	1468
8	Gross motor*.mp.	599
9	Psychomotor disorders/	1067
10	Developmental disabilities/ or motor skills disorders/	947
11	Developmental coordination disorder*.mp.	219
12	DCD.mp.	113
13	6 or 7 or 8 or 9 or 10 or 11 or 12	4982
14	5 and 13	1510
15	Clinical assessment scales/	4318
16	Questionnaires/	4123
17	Disability evaluation/	7023
18	Outcome measure*.mp/	9845
19	Outcome*.mp.	38379
20	Assess*.mp.	43680
21	Scale*.mp.	17562
22	Evaluat*.mp.	40621
23	15 or 16 or 17 or 18 or 19 or 20 or 21 or 22	93570
24	14 and 23	865
25	Measurement/	1629
26	Reproducibility of results.mp/	2241
27	"Consistency and reliability"/	1898
28	Statistics/	1075
29	Specificity.mp.	1241
30	Sensitivity.mp.	2860
31	"Predictive value of tests"/	839
32	25 or 26 or 27 or 28 or 29 or 30 or 31	10256
33	24 and 32	81

Supplementary table 5: Definition of terms

	Measurement Property	Definition	Example/explanation
Validity	Content	The degree to which an assessment tool's content measures the construct that it intends to measure ⁷	Concerned with the relevance and comprehensiveness of the items included in the assessment tool
	Construct	Measures the degree to which the scores obtained from the test are an adequate reflection of the construct to be measured ⁷	Examples include structural validity (whether scores reflect the dimensionality of the construct), hypothesis testing (item construct validity) and cross-cultural validity (whether translated or culturally adapted assessments adequately reflect the original version) ⁷
	Criterion	Assesses whether or not the test scores reflect a 'gold standard' assessment ⁷	As there is no gold standard of assessment for gross motor function in children this is often assessed with correlations of scores obtained from two or three other frequently used tools.
Reliability	Reliability	Refers to the consistency of a test score regardless of the time between assessments (test-retest) or the person administering (intra and inter-rater) ⁵⁰	Usually measured with intraclass correlation coefficient (ICC), but can be measured using Cohen's kappa coefficient. Percentage agreement and Pearson's correlation coefficient do not incorporate error into the calculations and as such is not a true measure of agreement ⁵⁰ . Scores > 0.80 are considered excellent, 0.60-0.79 adequate and <0.59 poor ¹¹
	Internal consistency	The degree of interrelatedness of an assessment tool's items ⁷	Usually measured using Cronbach's alpha (α) ⁷ . scores > 0.70 demonstrates high relationship, 0.5 to 0.69 a moderate relationship, 0.26 to 0.49 a low relationship and < 0.26 little relationship ⁵⁰ .
	Measurement Error	Refers to the error obtained between measurements that cannot be attributed to the patients true change ⁷	May be systematic or random error ⁷
Responsiveness	Responsiveness	An assessment tool's ability to detect change over time in the construct it purports to measure ⁷	This is central to a tools capacity to be used as an outcome measure.

Supplementary table 6: Excluded Assessment Tools

Reason	Assessments
Manual not available in English	Maastricht's Motor Test (MMT) The Motor-Proficiency-Test for children between 4 and 6 years of age (MOT 4-6) Zuk Assessment Körperkoordinationstest für Kinder (KTK)
Cannot extract meaningful gross motor score	Early Intervention Developmental Profile (EIDP) Neurological Developmental Exam Preschooler Gross Motor Quality Scale (PGMQ) The Malawi Developmental Assessment Tool (MDAT) Dutch table tennis motor skills assessment
Screening Tool	Brief Assessment of Motor Function (BAMF) The Motor Performance Checklist Motor skill checklist (MSC)
Diagnosis specific/requires a diagnosis	Assessment Battery for the Atypical Handicapped Child (VAB) Video-based documentation and rating system of the motor behaviour of handicapped children
Only assesses one motor domain (e.g. gait)	Standardized Walking Obstacle Course (SWOC) Timed floor to stand test
Manual not published/commercially available	Rapid Neurodevelopmental Assessment (RNDA) Tufts Assessment of Motor Performance (TAMP) Zurich Neuromotor Assessment (ZNA)

Supplementary table 7: Scoring and administration of assessment tools

Assessment Tool	Scoring	Interpretation of scores	Other
Bayley-III ¹²	Motor score - gross (varying items) and fine motor (varying items) subscales. Binary score with reverse/discontinue rules	Raw scores Composite scores Centile ranks Age equivalents Growth scores	Lends itself to multidisciplinary team testing.
BOT-2 ¹³	Fine manual (15 items) manual coordination (12 items) body coordination (16 items) strength and agility (10 items) subscales. Scoring differs for subtests	Raw scores Age adjusted standard scores Composite scores Centile ranks Age equivalents Descriptive categories. Complex conversions	Administration Easel includes instructions, diagrams and photos of test procedure
MABC-2 ¹⁴	Manual dexterity (3 items), aiming & catching (2 items) and balance (3 items) subscales.	Raw scores component scores centile ranks total test score traffic light system. Simple conversion	Also Available: MABC-2 Checklist (screening tool) and intervention manual
MAND ¹⁵	Fine motor (5 items) Gross motor (5 items)	Raw scores Scaled scores converted to an NDI. Factor scores. Complex conversions	Case studies included in manual for hyperactivity, encephalitis, mild head trauma, CP and muscular dystrophy
NSMDA ¹⁶	Functional grade given for each subscale, which is combined to create an overall score.	Indicates: normal range, minimal dysfunction, mild problems, moderate, severe or profound disability	Sections for comment on strengths, behavioural state during testing, musculoskeletal system and recommendations.
PDMS-2 ¹⁷	GM: Stationary (30 items), locomotion (89 items), object manipulation (24 item). FM: grasping(26 items) , visual-motor integration (72 items)	Raw scores, Age equivalent, centile rank. Standard scores (subtests) Composite quotient. Complex conversions.	Motor activities program (intervention ideas)
TGMD-2 ¹⁸	Locomotion (6 items) and Object Control (6 items). Separate male/female norms for object control subset	Raw scores, standard scores, percentile rank, age equivalent, Gross Motor Quotient. Simple conversion.	Simple to administer

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition ¹²; BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition ¹³; MABC-2, Movement Assessment Battery for Children 2nd edition ¹⁴; MAND, McCarron Assessment of Neuromuscular Development ¹⁵; NSMDA, Neurological Sensory Motor Developmental Assessment ¹⁶; PDMS-2, Peabody Developmental Motor Scales 2nd edition ¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition ¹⁸; GM, Gross Motor; FM, Fine Motor; NDI, Neurodevelopmental Index

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only



PRISMA 2009 Checklist

Section/topic	#	Checklist item	Reported on page #
TITLE			
Title	1	Identify the report as a systematic review, meta-analysis, or both.	1 Title page
ABSTRACT			
Structured summary	2	Provide a structured summary including, as applicable: background; objectives; data sources; study eligibility criteria, participants, and interventions; study appraisal and synthesis methods; results; limitations; conclusions and implications of key findings; systematic review registration number.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known.	4-5
Objectives	4	Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS).	5
METHODS			
Protocol and registration	5	Indicate if a review protocol exists, if and where it can be accessed (e.g., Web address), and, if available, provide registration information including registration number.	-
Eligibility criteria	6	Specify study characteristics (e.g., PICOS, length of follow-up) and report characteristics (e.g., years considered, language, publication status) used as criteria for eligibility, giving rationale.	5-6
Information sources	7	Describe all information sources (e.g., databases with dates of coverage, contact with study authors to identify additional studies) in the search and date last searched.	5
Search	8	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated.	5 and supplementary tables 2
Study selection	9	State the process for selecting studies (i.e., screening, eligibility, included in systematic review, and, if applicable, included in the meta-analysis).	5-6
Data collection process	10	Describe method of data extraction from reports (e.g., piloted forms, independently, in duplicate) and any processes for obtaining and confirming data from investigators.	6
Data items	11	List and define all variables for which data were sought (e.g., PICOS, funding sources) and any assumptions and simplifications made.	6-7
Risk of bias in individual studies	12	Describe methods used for assessing risk of bias of individual studies (including specification of whether this was done at the study or outcome level), and how this information is to be used in any data synthesis.	6
Summary measures	13	State the principal summary measures (e.g., risk ratio, difference in means).	7



PRISMA 2009 Checklist

Synthesis of results	14	Describe the methods of handling data and combining results of studies, if done, including measures of consistency (e.g., I^2) for each meta-analysis.	-
----------------------	----	---	---

Page 1 of 2

Section/topic	#	Checklist item	Reported on page #
Risk of bias across studies	15	Specify any assessment of risk of bias that may affect the cumulative evidence (e.g., publication bias, selective reporting within studies).	24, 25
Additional analyses	16	Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified.	-
RESULTS			
Study selection	17	Give numbers of studies screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally with a flow diagram.	Figure 1 + page 7
Study characteristics	18	For each study, present characteristics for which data were extracted (e.g., study size, PICOS, follow-up period) and provide the citations.	Table 1 – page 9 + Suppl table 3
Risk of bias within studies	19	Present data on risk of bias of each study and, if available, any outcome level assessment (see item 12).	11 + Table 3 – page 14-15
Results of individual studies	20	For all outcomes considered (benefits or harms), present, for each study: (a) simple summary data for each intervention group (b) effect estimates and confidence intervals, ideally with a forest plot.	7-8, 11-12, 17 + Table 2 page 13, Table 4 page 18, Table 5 page 19, Table 6 page 20 + Figures 2 & 3
Synthesis of results	21	Present results of each meta-analysis done, including confidence intervals and measures of consistency.	-
Risk of bias across studies	22	Present results of any assessment of risk of bias across studies (see Item 15).	-
Additional analysis	23	Give results of additional analyses, if done (e.g. sensitivity or subgroup analyses, meta-regression [see Item 16]).	-



PRISMA 2009 Checklist

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

DISCUSSION			
Summary of evidence	24	Summarize the main findings including the strength of evidence for each main outcome; consider their relevance to key groups (e.g., healthcare providers, users, and policy makers).	22-24
Limitations	25	Discuss limitations at study and outcome level (e.g., risk of bias), and at review-level (e.g., incomplete retrieval of identified research, reporting bias).	24
Conclusions	26	Provide a general interpretation of the results in the context of other evidence, and implications for future research.	25-26
FUNDING			
Funding	27	Describe sources of funding for the systematic review and other support (e.g., supply of data); role of funders for the systematic review.	1

From: Moher D, Liberati A, Tetzlaff J, Altman DG, The PRISMA Group (2009). Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA Statement. PLoS Med 6(7): e1000097. doi:10.1371/journal.pmed1000097

For more information, visit: www.prisma-statement.org.

BMJ Open

Psychometric properties of gross motor assessment tools for children: a systematic review

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2018-021734.R2
Article Type:	Research
Date Submitted by the Author:	23-Aug-2018
Complete List of Authors:	Griffiths, Alison; Monash University, Physiotherapy; The Royal Children's Hospital, Physiotherapy Toovey, Rachel; Murdoch Childrens Research Institute, Developmental Disability and Rehabilitation Research; University of Melbourne, Physiotherapy Morgan, Prue; Monash University, Physiotherapy Spittle, Alicia; University of Melbourne, Physiotherapy; Murdoch Childrens Research Institute, Victorian Infant Brain Studies
Primary Subject Heading:	Paediatrics
Secondary Subject Heading:	Rehabilitation medicine
Keywords:	PAEDIATRICS, REHABILITATION MEDICINE, validity, gross motor assessment, reliability

SCHOLARONE™
Manuscripts

Only

Psychometric properties of gross motor assessment tools for children: a systematic review

Alison Griffiths^{1,2,3}, Rachel Toovey^{3,4}, Prue E. Morgan¹, Alicia J. Spittle^{3,4}

1. Department of Physiotherapy, School of Primary and Allied Health Care, Monash University, Frankston, Victoria, Australia

2. Department of Physiotherapy, The Royal Children's Hospital, Parkville, Victoria, Australia

3. Murdoch Children's Research Institute, Parkville, Victoria, Australia

4. Department of Physiotherapy, The University of Melbourne, Parkville, Victoria, Australia

Corresponding Author:

Name: A/Prof Alicia J. Spittle

Address: Department of Physiotherapy, The University of Melbourne, Level 7 Alan Gilbert Building, 161 Barry Street, Parkville, Vic 3010, Australia

Email: aspittle@unimelb.edu.au

Funding: This study was part-funded by grants from the National Health and Medical Research Council Career Development Fellowship (AJS) 1053767 and Centre of Research Excellence in Newborn Medicine 1060733 (AJS and AG) and the Victorian Government's Operational Infrastructure Support Program.

Financial disclosure: The authors have declared that they have no financial relationships relevant to this article.

Conflict of interest: The authors have no conflict of interest.

Keywords: paediatrics, reliability, validity, rehabilitation medicine, gross motor assessment

Abstract

Objective:

Gross motor assessment tools have a critical role in identifying, diagnosing and evaluating motor difficulties in childhood. The objective of this review was to systematically evaluate the psychometric properties and clinical utility of gross motor assessment tools for children 2-12 years.

Method:

A systematic search of MEDLINE, Embase, CINAHL and AMED was performed between May and July 2017. Methodological quality was assessed with the COnsensus-based Standards for the selection of health status Measurement INstruments (COSMIN) checklist and an outcome measures rating form was used to evaluate reliability, validity and clinical utility of assessment tools.

Results:

Seven assessment tools from 37 studies/manuals met the inclusion criteria: Bayley Scale of Infant and Toddler Development-III (Bayley-III), Bruininks-Oseretsky Test of Motor Proficiency-2 (BOT-2), Movement Assessment Battery for Children-2 (MABC-2), McCarron Assessment of Neuromuscular Development (MAND), Neurological Sensory Motor Developmental Assessment (NSMDA), Peabody Developmental Motor Scales-2 (PDMS-2) and Test of Gross Motor Development-2 (TGMD-2). Methodological quality varied from poor to excellent. Validity and internal consistency varied from fair to excellent (α 0.5-0.99). The Bayley-III, NSMDA and MABC-2 have evidence of predictive validity. Test re-test reliability is excellent in the BOT-2 (ICC=0.80-0.99), PDMS-2 (ICC=0.97), MABC-2 (ICC=0.83-0.96) and TGMD-2 (ICC=0.81-0.92). TGMD-2 has the highest interrater (ICC 0.88-0.93) and intrarater reliability (ICC=0.92-0.99).

Conclusions:

The majority of gross motor assessments for children have good-excellent validity. Test-retest reliability is highest in the BOT-2, MABC-2, PDMS-2 and TGMD-2. The Bayley-III has the best predictive validity at 2 years of age for later motor outcome. None of the

1 assessment tools demonstrate good evaluative validity. Further research on evaluative gross
2 motor assessment tools are urgently needed.

3 Strengths and limitations of this study

- 4 • This systematic review comprehensively assesses methodological quality of included
5 studies using the COSMIN checklist.
- 6 • Results of this systematic review can provide guidance to clinicians when choosing
7 gross motor assessment tools based on test psychometric properties and clinical
8 utility.
- 9 • Areas for future research are identified including improving the evidence of inter and
10 intrarater reliability and responsiveness to change as well as the ascertainment of
11 predictive validity over a longer period of time.
- 12 • Only articles or test manuals written in English were included.
- 13 • Only one reviewer screened titles and abstracts for inclusion

1 Introduction

2 Motor function promotes cognitive and perceptual development in children and contributes
3 to their ability to participate in their home, school and community environments¹. Motor
4 impairment can negatively affect activity and participation levels of children², which may
5 lead to lower levels of physical activity, fitness and health into adulthood³. While severe
6 motor deficits are usually diagnosed before 2 years of age, mild motor deficits may not
7 become evident until children are in preschool and primary school environments where
8 they are exposed to increasingly complex tasks and compared to their peers³. Identification
9 of motor difficulties is an important step towards support and intervention for the child and
10 their family.

11 Healthcare professionals and researchers require standardised assessment tools to identify,
12 classify and diagnose motor problems in children⁴. Further, assessment tools are essential
13 to monitor the effects of interventions⁴. There is no gold standard of motor assessment for
14 children and the available tests vary in their ease of use and interpretability in clinical and
15 research settings, and whether they are norm or criterion referenced⁵. Criterion referenced
16 tests are designed to be scored as items or criteria are demonstrated; meaning that the
17 score is a reflection of a child's competence on the test items. Most available assessments
18 however, are norm referenced, meaning that a child's results are reported in relation to a
19 specific population⁴. The characteristics of the normed population should be taken into
20 consideration when interpreting test results as environmental and cultural differences have
21 been found to affect motor development⁶.

22 Health professionals should be aware of the validity and reliability of assessment tools to
23 assist in their instrument selection and interpretation of results. Validity refers to "The
24 degree to which [an instrument] is an adequate reflection of the construct to be measured"
25⁷. If an instrument does not have adequate construct or content validity then it may not be
26 assessing the skills that it purports to. Reliability refers to "the degree to which the
27 measurement is free from measurement error"⁷, which is significant when interpreting
28 results. If a child is assessed as being significantly delayed in their gross motor skills, the
29 reliability of that tool indicates the likelihood that a result is due to error.

1
2
3 1 A systematic review in 2010 by Slater⁸ evaluated performance-based gross motor tests for
4 2 children with developmental coordination disorder, however it did not include the second
5 3 and most recent version of the Movement Assessment Battery for Children 2 (MABC-2),
6 4 which is widely used. Brown and Lalor⁹ suggested that as a result of the changes to the
7 5 original Movement Assessment Battery for Children (MABC) in age range, age bands,
8 6 materials and tasks, that the MABC-2 requires independent reliability and validity
9 7 assessment. Over the past eight years there has also been a significant increase in the
10 8 number of papers assessing the psychometric properties of motor assessment tools in
11 9 children. A systematic review of these and previous papers is warranted, in order to add to
12 10 our understanding of the psychometrics of standardised gross motor assessment tools.

13
14
15
16
17
18
19
20
21 11 The primary aim of this systematic review is to identify and evaluate the clinical utility and
22 12 psychometric properties of gross motor assessment tools appropriate for use in preschool
23 13 and school age children from 2-12 years by assessing the methodological quality of the
24 14 included studies. The secondary aim of this review is to identify any areas for further
25 15 research.

16 Method

17 A comprehensive search strategy was completed in databases OVID Medline (1996 to May
18 2017), CINAHL plus (1937 to July 2017), Embase (1974 – May 2017) and AMED (1985 – July
19 2017) (Supplementary tables 1-4). The search strategy used MeSH terms and text words for
20 ('child' or 'paediatric') and ('motor skills' or 'motor activity' or 'gross motor' or
21 'psychomotor' or 'developmental coordination disorder') and ('questionnaires' or 'outcome
22 assessment' or 'instrument' or 'task performance') and ('reliability' or 'validity' or
23 'psychometrics'). Reference lists of included articles were also screened to identify any
24 additional papers. If full texts were unavailable or further information required regarding
25 availability of manuals authors were contacted.

26 Assessment tools were included if they were 1. Discriminative, predictive or evaluative of
27 gross motor skills, 2. Assessed \geq two gross motor (e.g. balance, jumping etc.) items, 3. Able
28 to extract a meaningful gross motor sub-score, 4. Applicable to children 2-12 years of age, 5.

1
2
3 1 Criterion or norm referenced test with a standardised assessment procedure and 6.
4 2 Instructional manuals are published or commercially available.
5
6
7 3 Articles describing use of the assessment tool were included if; $\geq 90\%$ of the study
8 4 population were within 2-12 years of age, it was available in English and if validity and/or
9 5 reliability of the assessment tool was reported.
10
11
12
13 6 Assessment tools were excluded if they met any of the following criteria 1. Questionnaires
14 7 or screening tools, 2. Only applicable to children with a specific diagnosis (e.g. cerebral
15 8 palsy, Down's syndrome), 3. Test manuals not available in English and 4. The version of the
16 9 test has been superseded.
17
18
19
20
21 10 Titles and abstracts were screened by the first author with any studies that clearly did not
22 11 meet inclusion criteria excluded. The remaining papers were obtained in full text and
23 12 reviewed by two authors (AG, RT or PM) with selection based on inclusion and exclusion
24 13 criteria. Papers and assessment tools were included after agreement by both raters, with
25 14 conflicting decisions discussed until a consensus was reached.
26
27
28
29
30 15 Methodological assessment of the papers was completed using the four-point scale of the
31 16 COnsensus-based Standards for the selection of health status Measurement INstruments
32 17 (COSMIN) checklist¹⁰. The COSMIN incorporates three quality domains: Validity, Reliability
33 18 and Responsiveness consisting of seven measurement properties: content, construct and
34 19 criterion validity, internal consistency, reliability, measurement error and responsiveness⁷
35 20 (Supplementary Table 5). Cross-cultural validity, structural validity and hypothesis testing
36 21 are all considered to be a component of construct validity⁷. Whilst predictive validity is
37 22 considered to be a component of content validity, it is reported on separately in this paper
38 23 for interpretability of results⁷.
39
40
41
42
43
44
45
46 24 The overall score for each measurement property on the COSMIN checklist is determined by
47 25 a 'worse score counts' approach¹⁰. Each property is rated as excellent, good, fair or poor
48 26 methodological quality based on descriptive criteria. Data extraction and assessment of
49 27 methodological quality was performed independently by two assessors (AG and RT). In the
50 28 case of any uncertainty a third reviewer (AS) performed a COSMIN assessment and
51 29 disagreement was resolved through discussion.
52
53
54
55
56
57
58
59
60

1
2
3 1 A data extraction form for each assessment tool was adapted from the CanChild Outcome
4 2 Measures Rating Form to collate information on clinical utility, validity, reliability and
5 3 responsiveness¹¹. Items chosen to represent the clinical utility of the assessment tools were
6 4 the cost of manuals, kits, training requirements, time to administer the assessment and the
7 5 ease of scoring. All reported values for reliability were collected, however, only those papers
8 6 reporting intraclass Correlation Coefficient (ICC) were directly compared.

13
14 7 Patient and Public Involvement

15
16 8 As this was a systematic review of existing papers there was no patients or public involvement.

19 20 21 9 Results

22
23 10 Figure 1 provides details of study selection. Seven assessment tools were identified for
24 11 inclusion; Bayley Scale of Infant and Toddler Development III (Bayley-III), Bruininks-
25 12 Oseretsky Test of Motor Proficiency 2 (BOT-2), Movement Assessment Battery for Children 2
26 13 (MABC-2), McCarron Assessment of Neuromuscular Development (MAND), Neurological
27 14 Sensory Motor Developmental Assessment (NSMDA), Peabody Developmental Motor Scales
28 15 2 (PDMS-2), and Test of Gross Motor Development 2 (TGMD-2). The corresponding manuals
29 16 were then added to the final yield resulting in thirty papers and seven manuals. Twenty
30 17 assessment tools were excluded (Supplementary Table 6).

31
32
33
34
35
36
37
38 18 The majority of assessment tools identified in this review are discriminative and most lend
39 19 themselves towards use in a research setting. All norm referenced tools are from western
40 20 countries and each identified test covers a different age range as shown in Table 1.

41
42
43
44 21 The TGMD-2 is the only tool that assesses gross motor skills in isolation and that focusses on
45 22 quality of performance. The other gross motor assessments were either in conjunction with
46 23 assessment of fine motor and/or balance (MAND, MABC-2, BOT-2 and PDMS-2) or as a
47 24 component of a developmental assessment (NSMDA, Bayley-III).

48
49
50
51 25 Despite the variability in test structures, there is some consistency of items included within
52 26 the gross motor skill subsets between tests. Most include a locomotion task such as walking,
53 27 running or stair climbing; an object control or manipulation task such as throwing or

1
2
3 1 catching a ball; and a static or dynamic balance task such as standing on one leg or hopping.
4 2 The PDMS-2, BOT-2 and the MAND also include strength assessments (the PDMS-2 only in
5 3 some age groups).
6
7
8 4 The number of gross motor items for assessment vary both within and between the tools
9 5 (Table 1). For example, the number of items tested in the Bayley-III and the PDMS-2
10 6 depends on the age and ability of the child. Several assessments report criteria for
11 7 describing gross motor delay, although all test manuals warn against diagnosing delay based
12 8 on a single assessment.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 1. Gross Motor Assessment Tool Characteristics

Assessment Tool	Domains Tested	Gross motor components tested	Age range	Diagnostic criteria	Primary purpose	Secondary purpose	Type of test	Normative sample (year)
Bayley-III ¹²	Gross motor, fine motor, cognitive, communication, social/emotional, adaptive	Static postures, dynamic movement, balance	1 mth – 3 yrs	Developmental delay: <25th centile or below 2SD. *	Discriminative	Predictive, Evaluative, Research tool	Norm	1700 children from the USA (2000)
BOT-2 ¹³	Gross motor, fine motor	Coordination, balance, running speed and agility, strength	4 – 21 yrs	*	Discriminative Evaluative	Research tool	Norm	1520 children from the USA (2005)
MABC-2 ¹⁴	Gross motor, fine motor, balance	Aiming and catching, static and dynamic balance	3 – 16 yrs	Traffic light system: Green = normal, amber = 'at risk' and red = definite motor impairment (<15%). *	Discriminative Evaluative	Intervention planning, Research tool	Norm	1172 children from United Kingdom (2006)
MAND ¹⁵	Gross and fine motor	Coordination, jumping, static and dynamic balance	3 yrs – 25 yrs	NDI 70-85 = mild 55-69 = moderate <55 = severe disability *	Evaluative	Research tool	Norm	2000 3-35 yrs from the USA (1970's)
NSMDA ¹⁶	Gross Motor, Fine Motor, Neurological, Postural Development, Infant Patterns of Movement, Sensory Motor. †	Sitting, kneeling, walking, balance, running, hopping, jumping, catching, motor planning	1 mth – 6 yrs	Total score 6-8 normal, 9-11 minimal, 12-14 mild, 15-19 moderate, 20-25 severe, >25 profound disability *	Evaluative Discriminative	Predictive, Research tool	Criterion	N/A
PDMS-2 ¹⁷	Gross motor, fine motor	Stationary (standing balance, sit-ups, push-ups), locomotion (walking, running, jumping, hopping, etc.), object manipulation (kick, throw, hit, catch)	Birth – 5 yrs	*	Discriminative Evaluative	Predictive, Research tool	Norm	2003 USA and Canada (1997-8)

1									
2									
3									
4									
5	TGMD-2 ¹⁸	Gross Motor	Locomotion (run, gallop, hop, leap, jump, slide) and Object control (batting, dribbling, catch, kick, throw, roll)	3 – 10 yrs	*	Discriminative Evaluative	Outcome measure, research tool, intervention planning	Norm	1208 USA children (1997-1998)
6									
7									
8									
9									

10 Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for
 11 Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental
 12 Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ NDI, Neurodevelopmental Index; SD, Standard Deviation; mth, month; yrs, years *, Advisable to use
 13 clinical reasoning; †, requires some manual handling; USA, United States of America

1
2
3 1 The PDMS-2 is notable for the inclusion of credit towards incomplete skills in the scoring system.
4 2 Most other tests award a point or credit towards a skill only if it is demonstrated to the full
5 3 satisfaction of the stated criteria (score of 0 or 1). The PDMS-2 however is scored 0-2 allowing for 1
6 7 mark to be allocated as a child progresses towards a skill without mastering it. The TGMD-2 is also
7 8 notable for its marking system, in which points are awarded for the quality of the action performed,
8 9 instead of satisfactory completion of the task only. These actions include preparatory movements
9 10 prior to running and jumping, or arm position during movements. The NSMDA marking criteria is
10 11 somewhat more complicated with a system of scores 1-4 with a symbol of “+” denoting hyperactive
11 12 response and “-” a hyporeactive response. The PDMS-2, MABC-2, BOT-2, MAND, TGMD-2 and
12 13 Bayley-III all require raw scores to be converted to a standard (or scaled) score based on tables
13 14 supplied in the manuals. For the BOT-2 this is a multiple step process which can then be converted to
14 15 both sex-specific or combined standard scores and percentile ranks. A summary of assessment tool
15 16 characteristics can be found in Table 1.
16 17
17 18
18 19
19 20
20 21
21 22
22 23
23 24
24 25
25 26

26 14 Clinical Utility

27 15 The clinical utility of the assessment tools is summarised in Table 2, while scoring and administration
28 16 is detailed in Supplementary Table 7. The shortest administration time is 15-20 minutes for the
29 17 TGMD-2 and the MAND; whilst most manuals report 20-60 minutes is required to complete an
30 18 assessment. These times are not inclusive of equipment set up, pack up and scoring, which varies
31 19 depending on the amount of equipment and complexity of the scoring process. All assessments
32 20 require the user to be familiar with the test before administration and to possess a high level of
33 21 understanding of child movement and development. The MABC-2 and PDMS-2 are the only
34 22 assessments that come with supporting material to guide intervention post assessment (when the
35 23 complete kit is purchased).
36 24
37 25
38 26
39 27
40 28
41 29
42 30
43 31
44 32
45 33
46 34
47 35
48 36
49 37
50 38
51 39
52 40
53 41
54 42
55 43
56 44
57 45
58 46
59 47
60 48

54 24 Methodological quality

55 25 All articles were assessed using the COSMIN checklist to determine methodological quality. Several
56 26 studies were marked down for failing to report missing data, small sample sizes and for using
57 27 inappropriate statistical methods. A summary of the articles and corresponding COSMIN
58 28 methodology rating is provided in Table 3.
59 29
60 30

Validity

The content and construct validity of the included assessment tools are summarised in Table 4.

Most assessments were developed by or with input from experts in the field, with most also performing literature reviews. Bruininks and Bruininks¹³ performed comprehensive surveys, pilot, tryout and standardisation studies before finalising the BOT-2, providing the most comprehensively reported content validity.

Construct validity was confirmed with factor analysis (either exploratory or confirmatory) in most assessment tools. The TGMD-2 has the most evidence for construct validity with several papers performing confirmatory and exploratory factor analysis^{19 20 18 21 22 23}. The MABC-2, BOT-2, Bayley-III, MAND and PDMS-2 had factor analysis performed only in one paper. The MABC-2 was shown to require changes to remain valid in the Chinese and Dutch speaking populations^{24 25}. The BOT-2, MABC-2 and TGMD-2 all provide evidence of the ability to discriminate between particular age or diagnosis groups, which can be considered to support their content validity. The NSMDA has minimal assessment of construct validity in children over 2 years. The Bayley-III, NSMDA and MABC-2 are the only assessments that provide evidence of predictive validity (Table 5). Concurrent validity between the MABC-2, PDMS-2 and BOT-2 is moderate to high, whilst the TGMD-2 is only weakly correlated with the MABC-2⁵ (Table 5). The PDMS-2, TMGD-2 and NSMDA report correlations with other criteria such as paediatrician diagnosis, physical fitness or psychomotor/intelligence tests.

Table 2. Clinical Utility of Gross Motor Assessment Tools

Assessment Tool	Time to administer (min)	Test Procedure	Target Examiner population	Training	Equipment/Manual
Bayley-III ¹²	30-90	Therapist administers in standardised order	Paediatric health professionals early childhood specialists	Formal training not required. DVD, webinars and workshops available	Comprehensive manual/kit: £1089 Test kit provides most equipment
BOT-2 ¹³	40-60	Therapist administered in standardised order	Paediatric health professionals early childhood specialists	Formal training not required	Comprehensive manual/kit: £961 Test kit provides most equipment
MABC-2 ¹⁴	20-40	Therapist administers items in standardised order. Some flexibility allowed.	Research psychologists, OT, PT, Paediatricians	Formal training not required.	Comprehensive manual/ kit: £1191 Test kit provides most equipment
MAND ¹⁵	15-20	Therapist administers items in standardised order.	Professionals e.g. education, neurology, OT, PT, psychology etc.	Formal training not required.	Manual and test kit: £1366 includes equipment
NSMDA ¹⁶	20-45	Observation followed by therapist administration of test items.	PT, OT	Formal training not required (but is available)	Comprehensive manual: £35. Equipment not included
PDMS-2 ¹⁷	45-60 (20-30 for GM only)	Standardised procedure.	Paediatric health professionals, PE teachers, early intervention specialists	Formal training not required	Comprehensive manual/kit: £553 Includes some but not all equipment required
TGMD-2 ¹⁸	15-20	Standardised procedure.	Teachers, health professionals (OT, PT, doctors)	Formal training not required	Kit includes manual and record form: £128. Equipment not included

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition¹²; BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition¹³; MABC-2, Movement Assessment Battery for Children 2nd edition¹⁴; MAND, McCarron Assessment of Neuromuscular Development¹⁵; NSMDA, Neurological Sensory Motor Developmental Assessment¹⁶; PDMS-2, Peabody Developmental Motor Scales 2nd edition¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition¹⁸; GM, Gross motor; OT, Occupational Therapy; PT, Physiotherapy; PE, Physical Education

Table 3. Methodological quality of included articles

Test	First author, Year	Country	Population (Age, Diagnosis)	Internal consistency	Reliability	Measurement error	Content validity	Structural validity	Hypothesis testing	Cross- cultural validity	Criterion validity	Responsive - ness
BAYLEY III	Bayley ¹²	USA	1-42 mths	Fair	Fair	Good	Excellent	Good	Good	-	Good	-
	Spittle, et al. ⁴	Australia	2,4 yrs, Ex prem	-	-	-	-	-	-	-	Good	-
	Visser, et al. ²⁶	Netherlands	2.2-10.8 yrs, GDD, L.I.	-	-	-	Excellent	Poor	-	-	-	-
BOT-2	Wuang and Su ²⁷	Taiwan	4-12 yrs ID	Excellent	Excellent	Excellent	-	-	-	-	-	Fair
	Wuang, et al. ²⁸	Taiwan	3-6 yrs ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Bruininks and Bruininks ¹³	USA	4-21 yrs	Good	Fair (interrater) Fair (test-retest)	Good	Excellent	Good	-	-	Good	-
MABC-2 (AB 1)	Ellinoudis, et al. ²⁹	Greece	3-5.5 yrs	Excellent	Good	-	-	-	-	-	-	-
	Hua, et al. ²⁴	China	3-6 yrs	Excellent	Good	-	Excellent	Excellent	-	Poor	Excellent	-
	Logan, et al. ⁵	USA	3-6 yrs	-	-	-	-	-	Fair	-	Fair	-
	Smits-Engelsman, et al. ³⁰	Belgium	3-4 yrs	Poor	Poor	Poor	-	-	-	-	-	-
	Holm, et al. ³¹	Norway	7-9 yrs	-	Fair (interrater) Poor (intrarater)	Poor	-	-	-	-	-	-
MABC-2 (AB 2)	Kita, et al. ³²	Japan	7-10 yrs	Excellent	-	-	-	-	-	Poor	-	-
	Griffiths, et al. ³³	Australia	4-8 yrs	-	-	-	-	-	-	-	Good	-
MABC-2	Henderson, et al. ¹⁴	UK	3-16 yrs	-	Fair	Good	Excellent	-	-	-	-	-
	Niemeijer, et al. ²⁵	Netherlands + Belgium	-	-	-	-	-	-	-	Poor	-	-
	Schulz, et al. ³⁴	U.K	3-16 yrs	-	-	-	Excellent	Good	-	-	-	-
	Valentini, et al. ³⁵	Brazil	3-13 yrs	Fair	Fair	-	Fair	Poor	-	Poor	Poor	-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

	Wuang, et al. ²⁸	Taiwan	3-6 yrs, ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Wuang, et al. ³⁶	Taiwan	6-12 yrs DCD	Poor	Fair	Good	-	-	-	-	-	Fair
MAND	Hands, et al. ³⁷	Australia	10-17 yrs	-	-	-	-	Excellent	-	-	-	-
	McCarron ¹⁵	USA	7yrs	-	-	-	Fair	Poor	-	-	Poor	-
NSMDA	Danks, et al. ³⁸	Australia	2 + 4 yrs ELBW	-	-	-	-	-	-	-	Fair	-
	MacDonald and Burns ³⁹	Australia	2 + 4 yrs CP	-	-	-	-	Fair	-	-	Poor	-
	Burns, et al. ⁴⁰	Australia	1-24 mths VLBW	Poor	-	-	Poor	-	-	-	-	-
	Burns, et al. ⁴¹	Australia	1-mnths VLBW	-	-	-	-	Poor	-	-	Fair	-
PDMS-2	Hua, et al. ²⁴	China	3-6 yrs.	Excellent	Good	-	Excellent	Excellent	-	Poor	Excellent	-
	Wuang, et al. ²⁸	Taiwan	3-6 yrs ID	Fair	Good	Good	-	-	-	-	Good	Fair
	Folio and Fewell ¹⁷	USA	0-71 mths	Good	-	Poor	Excellent	Good	Good	-	Poor	-
TGMD-2	Barnett, et al. ⁴²	Australia	4-8 yrs	-	Fair	-	-	-	-	-	-	-
	Farrokhi, et al. ⁴³	Iran	3-11 yrs	Fair	Fair	-	Fair	Fair	-	-	-	-
	Houwen, et al. ²¹	Netherlands	6-12 yrs VI	Fair	Fair	-	-	Fair	-	-	-	-
	Kim, et al. ⁴⁴	Korea	8-12 yrs ID	-	Poor	-	-	-	-	-	-	-
	Kim, et al. ⁴⁵	Korea	5-6 yrs	Poor	Fair	-	-	Poor	-	-	Poor	-
	Logan, et al. ⁵	USA	3-6 yrs	-	-	-	-	-	Fair	-	Fair	-
	Rudd, et al. ¹⁹	Australia	6-12 yrs	-	-	-	-	Good	-	-	-	-
	Simons, et al. ²³	Belgium	7-10 yrs ID	Good	Good (interrater) Poor (test-retest)	-	Excellent	Good	Good	-	-	-
	Valentini ²⁰	Brazil	3-10 yrs	Poor	Fair (test-retest) Good (intra,	-	Excellent	Good	-	Fair	Good	-

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

interrater)											
Wong and Yin Cheung ²²	China	3-10 yrs	-	-	-	-	Fair	-	-	-	-
Ulrich ¹⁸	USA	3-10 yrs	Good	Fair (test-retest)	Fair	Poor	Good	-	-	Fair	-
Poor (interrater)											

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ Mths, Months; yrs, years; DCD, Developmental Coordination Disorder; VI, Vision Impairment; ID, Intellectual Disability; GDD, global developmental delay; L.I, Language Impairment; ELBW, Extremely Low Birth Weight; VLBW, Very Low Birth Weight; CP, Cerebral Palsy; prem, premature; USA, United States of America

For peer review only

1 Reliability

2 Internal consistency of assessments are summarised in Table 6. The BOT-2's high internal
3 consistency is well supported, including for children with an intellectual disability^{28 46}. The MABC-2
4 appears to have lower internal consistency than the BOT-2, which may relate to the limited number
5 of test items (eight) on the MABC-2. The highest values for internal consistency for the MABC-2
6 were obtained in specific populations (intellectual disability and developmental coordination
7 disorder) with poor to fair methodology only. Conversely the highest quality articles reported the
8 lowest values, although it should be noted that these assessed age band 1 (3-6 years) only. Internal
9 consistency is reported to be high for the PDMS-2, while the Bayley-III is shown to have excellent
10 internal consistency in children aged 24-42 months.. The TGMD-2 is reported by two good quality
11 (and four poor to fair quality) articles to have excellent internal consistency, including for children
12 with vision impairment and intellectual disability. The MAND is the only assessment tool included in
13 this review without published data of internal consistency or reliability in this age group.

14 The reliability findings are summarised in Table 6 and in Figures 2 and 3. Test-retest reliability was
15 excellent in the Bayley-III (Table 6), BOT-2 and PDMS-2; and was good to excellent in the MABC-2
16 and TGMD-2 (Figure 2). Intra-rater reliability was rarely investigated or reported for most tools, with
17 the TGMD-2 demonstrating better results than the MABC-2 (Figure 3). Only the TGMD-2 and
18 MABC-2 report inter-rater reliability values using an ICC (Figure 3)^{31 42}. Inter-rater reliability is also
19 supported in the BOT-2 with Pearson Correlation Coefficient and Kappa respectively. The studies
20 referred to in the test manuals for the TGMD-2, Bayley-III, BOT-2 and MABC-2 all report reliability
21 findings using Pearson's correlation, which is less ideal than an ICC or weighted kappa for statistical
22 analysis^{47 48}. Only studies reporting ICC's are visually represented in Figures 2 (test-retest) and 3
23 (inter and intra-rater). The TGMD-2 test-retest reliability results from Houwen, et al.²¹ were
24 believed to contain an error as the reported ICC was outside of the reported confidence intervals
25 (ICC 0.92, 0.82-0.91). This data set was therefore excluded from Figure 2.

26 Responsiveness was reported for the Bayley-III, BOT-2, MABC-2 and PDMS-2 with minimal
27 detectable change (MDC) or a standard error of measurement (SEM)²⁸. Sensitivity and specificity
28 for detecting change was shown to be satisfactory in the MABC-2, PDMS-2 and MABC-2²⁸ (Table 6).
29 There have been no studies to date on the responsiveness of the TGMD-2, NSMDA or MAND.

Table 4: Content and construct validity of assessment tools

Test	Content	Construct
BAYLEY III	Expert opinion for standard and low verbal version ^{12 26} . Literature reviews. Gross motor score correlated with Motor component 0.70 ¹²	Factor analysis. Difference in mean scores with pervasive developmental disorder, and specific language impairment ¹² . H_i (gross motor subset) = 0.52-0.97 for children with language impairment and 0.82-0.99 in control group ²⁶
BOT-2	Focus groups, product survey, pilot, national tryout and standardisation studies, professional reviews ¹³	Factor analysis, scores increase with age, discriminates between normal and children with DCD ($N=50$), high-functioning ASD ($N=45$) and mild-moderate ID ($N=66$) ¹³
MABC-2	Expert Panel, Stakeholder feedback, Literature review ³¹ Expert panel - clarity (validity content index 71.8-93.9, Kappa 0.76-0.88) and pertinence (98.5-99.3 and kappa 0.83-0.92) $p<0.001$ ³⁵	Factor analysis, correlation coefficients ²⁹ Subtest correlations 0.65-0.76 $p<0.001$. Discriminates between ASD and control group ³¹ . Structural equation modelling (for each age group) ³⁴ . Expert panel - adequate face validity ³⁵ . Significant difference between TD, DCD and at risk DCD scores ($\eta^2 = 0.63$) $p<0.0001$ ³⁵ . UK norms not appropriate to use with Dutch/Flemish children as under/over-estimate risk of motor impairment ²⁵ . In Chinese population: CFA initially rejected. Acceptable fit achieved after 2 items removed ²⁴ . Age band 2 shows good validity in Japanese population ³² .
MAND	Based on neuropsychological theory. Several rounds of revision/trials of tasks during development ¹⁵	Factor analysis ^{15 37} . Scores increase with age, and discriminate between typically developing children and those with head trauma or neurological dysfunction as well as gender ^{15 37}
NSMDA	Literature review. Developed by an experienced paediatric physiotherapist ⁴⁰	Factor analysis (up to 2 years of age) ^{40 41} . Stability of test results over time (up to 2 years) ^{40 41} .
PDMS-2	Literature review. Created by experts in the field. Revised with feedback from therapists guided revision. Hierarchical sequence of items ¹⁷	Item response modelling. Factor analysis. Differential item functioning analysis. Scores correlated with age ($r=0.80-0.93$) ¹⁷
TGMD-2	Expert Panel (3 PE teachers with post-grad qualifications) ¹⁸ . Translated version (Brazilian Portuguese) language clarity 0.96, pertinence >0.89 . Experts CVI for clarity and pertinence were also strong- $\alpha = 0.93$ clarity and $\alpha = 0.91$ pertinence ²⁰	Exploratory and confirmatory factor analysis ^{19 20 18 21 22 23} High and significant correlation of increasing age and increasing scores ⁴³ . Age and disability differentiation ^{18 23} Subtest correlation 0.41 ¹⁸ Gallop, running and leaping not well correlated with locomotion subscale. Object control significant & highly correlated ⁴⁵ . ANOVA - significant age effect for object control ²³ Moderate correlation between items and subset scores, and between subset scores and total score ²³

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸; H_i , scalability coefficient; CFA, Confirmatory Factor Analysis; TD, Typically Developing; ASD, Autism Spectrum Disorder, ID, Intellectual Disability; WPPSI, *Wechsler Preschool and Primary Scale of Intelligence*; WISC-R, Wechsler Preschool and Primary Scale of Intelligence-R; NDI, Neurodevelopmental Index; ANOVA, Analysis of Variance

Table 5: Criterion and predictive validity of assessment tools

Test	Criterion	Predictive
BAYLEY III	Given but mean age <22 months. Not relevant to study population. ¹²	Motor impairment at 4 years: Bayley III at 2 years <1SD = sensitivity 0.32-0.037 specificity 0.97 <2SD sensitivity 0.18-0.21 specificity 1.00. CP at 4 years: Bayley III at 2 years <1SD sensitivity 0.83 specificity 0.94. <2SD sensitivity 0.67 specificity 1.0 ⁴
BOT-2	MABC-2 $\rho = 0.92$ PDMS-2 $\rho = 0.88$ ($N = 38$) ²⁸ . PDMS-2 Total motor composite $r = 0.77$ ¹³ .	-
MABC-2	PDMS-2 $\rho = 0.631 - 0.84$ ^{28,24} . TGMD-2 $\rho = 0.45$ ⁵ . TGMD-2 standard scores ($r = 0.3, p < 0.02$) ³⁵ . BOT-2 $\rho = 0.90 - 0.92$ ²⁸ .	Classification groups (DCD, at risk and TD) remained same over time (6 months) $\chi^2 = 0.67 p = 0.72$ ³⁵ . Predictive of motor impairment over 6-12 months ($N=41$) ICC 0.88 $p < 0.007$ ³⁵ . Scores at 4 years predictive of motor impairment at 8 years in children born <30 weeks gestation (PPV 79, sensitivity 79%, specificity 93%) ³³
MAND	Gross motor subscore: Low-moderate correlation with manual dexterity (-0.46 to 0.35), reaction time (-0.31 to -0.58), intelligence measures (WISC-R, Metropolitan Achievement Test) (0.30-0.39) and visual motor test (-0.33 to 0.39) ¹⁵	-
NSMDA	NSMDA at 2 years ($N = 148$) predictive of medical diagnosis $\chi^2 = 0.08 p = NS$ ⁴¹	Motor outcome at 11-13 yrs. NSMDA at 2years - sensitivity 48.8%, specificity 82.4%, NSMDA at 4 years sensitivity 64.5%, and specificity 80%. PPV at 2 years 83% at 4 years 87% ³⁸ . If classified 'severe' at 24 months - approximately 50% chance walking at 4 years (moderate = 80%, mild = 93% minimal = 100%) ³⁹
PDMS-2	MABC-2 $\rho = 0.63 - 0.84$, ^{24,28} MABC-2 gross motor composite $\rho = 0.743$ ²⁴ BOT-2 $\rho = 0.88$ ²⁸ . Mullen Scales of Early Learning GMQ = 0.86 FMQ = 0.80 ¹⁷	-
TGMD-2	MABC-2 total $r = 0.49 p < 0.01$ ⁵ . 'Teacher report' $r = 0.34-0.45$. physical fitness $r = -0.47 - 0.55$ ⁴⁵ ($N=41$) Basic Motor Generalizations subtest of the CSSA $r = 0.63$. Locomotor 0.63 object control 0.41 ¹⁸	-

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ NS, Not Specified; SD, Standard Deviation; CP, Cerebral Palsy; TD, Typically Developing; ICC, Intraclass Correlation Coefficient; χ^2 , Chi Squared; NDI, Neurodevelopmental Index; CSSA, Comprehensive Scales of Student Abilities

Table 6: Reliability of assessment tools

Test	Internal Consistency	Test-Retest	Intra-rater	Inter-rater	Minimal detectable change	Minimal clinical important difference
BAYLEY III	GM $\alpha = 0.87$ - 0.93 MC: $\alpha 0.90$ - 0.96 (24-42 months) ¹²	Gross Motor subtest (N=47) $r=0.79$ Motor component $r=0.80$ ¹²	-	-	SEM Gross motor subtest 0.85-1.08. of Motor component = 3.00-4.74 (24-42 months) ¹²	-
BOT-2	(N = 100) $\alpha = 0.92$ ²⁷ (N = 141) $\alpha = 0.86$ ²⁸ 4-7 yrs (N= 620) $\alpha = 0.95$ 8-11 yrs (N= 450) $\alpha = 0.95$ ¹³	(N = 100) ICC = 0.99 ²⁷ (N = 141) ICC = 0.97 ²⁸ 4-7 yrs (N = 43) $r = 0.81$ (8-12 yrs (N= 44) $r = 0.80$ ¹³	-	Total motor composite 4-21 yrs (N = 47) $r = 0.98$ ¹³	4.18 (sensitivity 55.10% specificity 72.55%) ²⁷ 7.43 (sensitivity 42.49% specificity 65.72%) ²⁸	6.53 (sensitivity 48.98% specificity 76.47%) ²⁷ 6.55 (sensitivity 49.99% specificity 58.78%) ²⁸
MABC-2 (AB 1)	(N = 60) M.D $\alpha = 0.51$, A&C $\alpha = 0.70$, Bal $\alpha = 0.66$ ²⁹ (N = 1823) $\alpha = 0.502$ ²⁴ (N=50) $\alpha = 0.81$ - 0.87 ³⁰	(N=60) ICC = 0.85 ²⁹ Item ICC's 0.830-0.985 ²⁴ ICC test-retest = 0.83 ³⁰ Inter-rater test-retest ICC = 0.79 ³⁰	(N=28) $\kappa = 0.71$ ³⁰	Item ICC's range 0.892-0.998 ²⁴ (N=22) $\kappa = 0.60$ ³⁰	(N=28) Intrarater MDC = 3.43 (N=22) Inter-tester MDC = 3.81 ³⁰	-
MABC-2 (AB 2)	Translated version (Japanese) (N=132) $\alpha = 0.602$ ³²	-	ICC = 0.64 ³¹	ICC 0.63 ³¹	Intra-rater SDC TTS: +/- 11.7 TSS +/- 3.3. Inter-rater SDC TTS +/-16.0 TSS +/- 3.8 ³¹	-
MABC-2	Subscales $\alpha = 0.78$ (M.D = 0.77, BS = 0.52, Bal = 0.77) ³⁵ $\alpha = 0.88$ ³⁶ (N = 141) $\alpha = 0.88$ ²⁸	N=60 (all 3 age bands) $r=0.80$ ¹⁴ $r=0.74$ $p<0.0001$ (standard score). ICC standard score = 0.85 ³⁵ ICC 0.96 ³⁶ N = 141 ICC =0.96 ²⁸	ICC 0.88 ³⁵	ICC 0.96-0.99 ³⁵	SEM 1.34 (95%CI) = 3 ¹⁴ 1.83 (95%CI) ³⁶ 1.83 (sensitivity 69.69% specificity 52.10%) ²⁸	1.39 (sensitivity 72.47% specificity 46.18%) ^{28 36}
MAND	-	-	-	-	-	-
NSMDA	Cross correlation matrix Item scoring (12+24months) 0.73 $p<0.001$, Functional grade (12+24months) 0.87 $p<0.001$ ⁴⁰	-	-	-	-	-
PDMS-2	(N=141) $\alpha=0.89$ ²⁸ 24-35m $\alpha=0.97$, 36-47m $\alpha=0.95$, 48-59m $\alpha=0.97$, 60-71m $\alpha=$	N=141 ICC= 0.97 ²⁸	unable to extract data for ≥ 24 months	unable to extract data for ≥ 24 months ¹⁷	7.76 (sensitivity 60.65% specificity 74.13%) ²⁸ SEM 24-	8.39 (sensitivity 61.65% specificity 71.34%) ²⁸

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

	0.98. For subgroups† $\alpha=0.99$ ¹⁷		¹⁷	59 months = 3, 60-71m = 2 ¹⁷
TGMD-2	($N=1438$) $\alpha=0.80$ ⁴³ $N=75$ Locomotor subset $\alpha=0.71$ object control $\alpha=0.72$ ²¹ $N=120$ $\alpha = 0.72$ ⁴⁵ $N=99$ $\alpha = 0.90$ ²³ $N = 1208$ Cronbach's $\alpha = 0.91$ (gross motor quotient). Locomotor 0.85 and object control 0.88. Note SEM GMQ = 4-5 SEM subsets=1 ¹⁸	$N=63$ ICC=0.81 95% CI ⁴³ $N=23$ ICC=0.92 total 95% CI ²¹ $N=99$ $r=0.98$ ²³ Locomotor test $r = 0.90$ $p < 0.0001$ object control test $r = 0.91$ $p < 0.001$ ²⁰ $N = 75$ $r=0.96$ overall (3-5 yrs $r = 0.91$), 6-8 years $r = 0.95$), (9-10 years $r = 0.94$) ¹⁸	$N=32$ ICC=0.97 95% CI ⁴³ $N=25$ ICC=0.95 95% CI ²¹ ICC = 0.78 ⁴⁴ ICC=0.92-0.99 ²⁰	Obj ICC=0.93 ⁴² ($N=50$) ICC=0.89 ²¹ ICC=0.75 ⁴⁴ $N=8$ $r= 1.00$ ²³ L.S ICC=0.88 Obj ICC=0.89 ²⁰ $N = 30$ $r=0.98$ ¹⁸

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition;¹² BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition;¹³ MABC-2, Movement Assessment Battery for Children 2nd edition;¹⁴ MAND, McCarron Assessment of Neuromuscular Development;¹⁵ NSMDA, Neurological Sensory Motor Developmental Assessment;¹⁶ PDMS-2, Peabody Developmental Motor Scales 2nd edition;¹⁷ TGMD-II, Test of Gross Motor Development 2nd edition;¹⁸ GM, Gross Motor Subset; MC, Motor Component; K, Kappa Coefficient; M.D, Manual Dexterity; BS, Ball Skills; BAL, Balance; A&C, Aiming and catching; SDC, Smallest Detectable Change; TTS, Total Test Score; TSS, Total Standard Score; †, gender, ethnicity, speech/language or physical disorder; Obj, Object Control Subset; L.S, Locomotion Subset

Discussion

This review identified seven gross motor assessment tools appropriate for use in clinical or research settings, each with their own strengths and limitations. Interestingly, only one of the seven assessments (TGMD-2) measured gross motor skills in isolation. This is likely a reflection on current practice to assess children's development as a whole, rather than assessing individual domains in isolation. A gross motor assessment embedded within a developmental assessment, such as that of the Bayley-III may be more appropriate than an isolated gross motor assessment for children where there is suspicion of multiple impairments.

A review by Slater, et al.⁸ reported that the TGMD-2 and the MABC (first edition) were recommended for assessing gross motor skills in children with developmental coordination disorder, but found that the MABC needed further evidence of validity. Cools, et al.⁴⁹ also published a detailed review of the clinical utility of gross motor assessment tools for children, but did not address the validity, reliability or responsiveness to change of these measures. This review adds to the literature by including updated information on the psychometric properties of the measures and a thorough methodological assessment using the COSMIN checklist which allows the reader to interpret these results with confidence. We have identified ten additional publications to support the content, construct and criterion validity of the MABC-2 and have demonstrated an overall higher methodological quality of the papers assessing the MABC-2 when compared with the TGMD-2. Papers that received lower methodological scores on the COSMIN can be attributed to inadequate reporting statistical methods, small sample sizes and non-independent assessors. Further research in this area should consider addressing these limitations in their study design to reduce potential error and increase confidence when interpreting results.

Content validity has been established for five of the included assessment tools, however, further research into the content validity for the MAND and NSMDA is required. The NSMDA's ability to predict a diagnosis of CP and motor outcomes over time does support its content validity, however the methodology scored as poor to fair on the COSMIN and as such content validity cannot be fully established. The use of expert panels, focus groups and/or stakeholder feedback for the BOT-2, MABC-2, TGMD-2 and PDMS-2 demonstrate thorough consideration of the relevance and comprehensiveness of the each test's assessment items during development.

1 The TGMD-2 is the only assessment tool considered to have well established construct validity, with
2 several papers reporting factor analysis. The NSMDA has undergone factor analysis for children up to, but
3 not beyond two years of age and as such further research is needed to support its validity in older children.
4 All other included assessment tools have undergone factor analysis assessment of their construct validity in
5 one paper and are supported by the ability to discriminate between medical diagnosis or age, and as such
6 are considered to have adequate construct validity. The criterion validity indicates that the TGMD-2 may
7 be measuring a slightly different construct to the other assessment tools included in this study as it has
8 poor agreement with the MABC-2, which in turn has good agreement with the PDMS-2 and the BOT-2.
9 This difference may be related to the inclusion of the assessment of quality of movement in the TGMD-2,
10 or the inclusion of balance and/or fine motor tasks on the other assessments. There is scope to investigate
11 the criterion validity of the MAND and the gross motor subsections of the Bayley-III and the NSMDA with
12 the other assessment tools in this study in the future.

13 The BOT-2 was the only assessment tool to have its reliability assessed with excellent methodology. In
14 conjunction with its reported results it can be considered to have the strongest evidence for internal
15 consistency and test-retest reliability out of the included assessment tools. The PDMS-2 and the MABC-2
16 can be considered to have the next best established test-retest reliability with good methodological
17 quality. The reported test-retest reliability values for the TGMD-2 are impacted by the poor to fair
18 methodological quality, and further high quality research needs to be done to support its body of evidence.
19 Test-rest, inter or intra-rater reliability has not been assessed in the MAND and NSMDA. In the clinical
20 context gross motor assessments are often repeated over time or between therapists and as such these
21 measures of reliability should be established. The Bayley-III would also benefit from further research into
22 its reliability, with no published inter or intra-rater reliability measures, and with only one, fair quality
23 report of good test-retest reliability.

24 As yet there is little evidence to support the use of these assessments as outcome measures. The inclusion
25 in some of the articles of minimal detectable change (MDC) and minimal clinically important difference
26 (MCID) is valuable for clinicians⁷. The difference between MDC and MCID is also of importance, as a change
27 in score does not necessarily relate to a meaningful change for the child or their family. Only the Bayley,
28 BOT-2, MABC-2 and PDMS-2 have a reported MCID with satisfactory sensitivity and specificity, however,
29 due to the fair methodological quality used to obtain these values they cannot be utilised with a high level
30 of confidence until further studies have been performed. The TGMD-2 was created in part to be used as an
31 outcome measure, however there are no articles to date investigating its responsiveness to change¹⁸. It
32 should also be noted that all of the included assessment tools measure impairment and activity limitations,
33 but do not specifically address the other elements of the International Classification of Functioning,

1 Disability and Health (ICF) domains of participation, personal factors and environment ². Clinicians should
2 utilise appropriate assessments or questionnaires to ensure that these domains of health are also
3 addressed in line with World Health Organisation guidelines ².
4

5
6 4 When considering a test's reliability all three elements of test error should be taken into account – these
7 can be described as time sampling (assessed with test-retest reliability), content sampling (assessed as
8 internal consistency), and inter-scorer difference (or interrater reliability) ¹⁸. This is one of the reasons that
9 clinicians should consider repeating assessments and/or completing a second alternative assessment. All
10 assessments should be interpreted in conjunction with clinical reasoning and observation. Included
11 assessment tools are not intended to be diagnostic on their own; results need to be combined with other
12 assessments and expert opinion to arrive at a clinical diagnosis.
13
14

15
16 11 All of the included assessment tools were found to have merits and limitations in their clinical utility the
17 body of evidence to support their use. Clinicians and researches should select their assessment tool with
18 consideration of psychometric properties (inclusive of the methodological rigour behind them), clinical
19 utility and for the population, situation and age group in question.
20
21

22
23 15 A potential limitation of this study was that one author screened the titles and abstracts, which may have
24 led to a sampling bias. Whilst care was taken to include all potentially relevant papers and assessment
25 tools until the second round of assessment with two authors, the potential for exclusion of papers relevant
26 to this review remains. The process of excluding both papers and assessment tools in this single step may
27 also be seen as a limitation, as the total number of assessment tools (or different versions of tools) was not
28 reported. This process does, however comply with the COSMIN and PRISMA guidelines. A second
29 limitation was the restriction of included papers and manuals to those published in English. Unfortunately
30 this resulted in the exclusion of three assessment tools that have been reported as commonly used in
31 Europe: The Motoriktest für Vier- bis Sechsjährige Kinder (MOT 4-6), the Körperkoordinationstest für Kinder
32 (KTK) and the Maastrichtse Motoriek Test (MMT)⁴⁹. The authors also note the third edition of the TGMD is
33 soon to be published and will need to be subjected to a similar level of assessment of psychometric
34 properties in the future.
35
36

37
38 27 Clinicians and parents who need guidance to set realistic therapy goals and to understand future
39 intervention requirements benefit from understanding a test's predictive ability. The NSMDA and the
40 MABC-2 are the only tools that have demonstrated long term (≥ 4 years follow up) predictive validity, while
41 the Bayley-III has good predictive validity at 2 years for future movement difficulties and for the diagnosis
42 of cerebral palsy at 4 years. However, further research into the long-term predictive validity of all included
43 gross motor assessment tools is warranted.
44
45

1 While validity and reliability should guide selection of assessment tools, clinical utility must also be taken
2 into consideration. Most tests have ongoing costs associated with forms and equipment replacement,
3 which may be prohibitive to some users. The NSMDA requires the therapist to handle the child for several
4 items which should be considered in relation to manual handling policies of institutions. Assessment
5 burden for children and families should also be taken into consideration when selecting an assessment
6 tools. Younger children are more likely to be distracted and may not understand test items as well, which
7 may also increase assessment times³⁰.

8 When a new edition of an assessment tools is released resulting in a change in age groups, scoring or tasks
9 it is insufficient to rely on the psychometric assessments that were performed on the original test. The
10 MABC-2 manual provides justification for the inclusion of reliability and validity assessment of the original
11 MABC¹⁴, however, owing to the significant changes in age groups and tasks between editions these were
12 not included for the analysis of the MABC-2 in this review. Two studies quoted in the MABC-2 manual to
13 support the validity and reliability are both unpublished works and as such are also unable to be included
14 in this systematic review. This could indicate a publication for the MABC-2.

15 The thorough methodological assessment of the included articles using the COSMIN checklist should be
16 seen as a strength of this paper, as should the range of assessment tools included in this review. While it
17 has previously been argued that the 'worst score counts' criteria in the COSMIN creates a floor effect⁵⁰,
18 the COSMIN authors argue that only 'fatal flaws' contribute to an overall score of poor¹⁰. There are few
19 tools available to assess the psychometric properties of assessment tools and arguably none so robustly
20 validated as the COSMIN.

21 There are many appropriate gross motor assessment tools available for use in research and clinical settings
22 today. Most of the available tools demonstrate adequate validity and reliability in children aged 2-12 and
23 as such the authors do not believe that new assessment tools need to be developed for use. There is scope
24 however to improve the evidence of inter and intra-rater reliability and predictive validity should be
25 ascertained over a longer period of time and with greater methodological rigour. Tools also need clearer
26 assessment of their responsiveness to change to assist clinicians and researchers with outcome measure
27 selection. Researchers should be mindful of the methods they use to assess validity and reliability. Clarity
28 of reporting, statistical methods and sample sizes should be carefully considered to ensure the highest
29 quality of evidence.

Conclusion

Currently available gross motor assessment tools for children have good to excellent content and construct validity. The BOT-2, MABC-2, PDMS-2 and TGMD-2 are the most reliable assessments in this age group. The Bayley-III has the best predictive validity at 2 years of age, and the NSMDA and the MABC-2 both have good predictive validity at 4 years of age. There is scope for further research into the predictive validity, reliability and responsiveness of gross motor assessment tools in preschool and school aged children. In practice clinicians should choose assessments with consideration of their psychometric properties in the context of the child that they are assessing.

Author Contributions

All individuals listed as authors meet the appropriate authorship criteria and have approved the acknowledgement of their contributions. The primary author, Ms Griffiths was responsible for the drafting of the paper and liaising with the co-authors on findings and conclusions. Ms Toovey contributed to the paper through interpretation of data, completing methodological assessments and revising manuscript content throughout its development. A/Profs Morgan and Spittle both contributed to the paper through assisting with the development of research design, interpretation of data and revising manuscript content through its development.

Data Sharing Statement

This paper includes data obtained from reviewing papers of published manuscripts. Data can be accessed by contacting the primary author.

Figures

Figure 1. PRISMA flow diagram detailing study selection

Figure 2. Test re-test reliability of gross motor assessment tools

Figure 2 legend: BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition¹³; MABC-2, Movement Assessment Battery for Children 2nd edition¹⁴; PDMS-2, Peabody Developmental Motor Scales 2nd edition¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition¹⁸.

- 1 1 Figure 3. Inter and interrater reliability of gross motor assessment tools
- 2
- 3 2 Figure 3 legend: MABC-2, Movement Assessment Battery for Children 2nd edition ¹⁴; TGMD-II, Test of Gross
- 4 3 Motor Development 2nd edition ¹⁸
- 5
- 6
- 7 4
- 8
- 9
- 10 5
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

For peer review only

References

1. Piek JP, Baynam GB, Barrett NC. The relationship between fine and gross motor ability, self-perceptions and self-worth in children and adolescents. *Human Movement Science* 2006;25(1):65-75. doi: <http://dx.doi.org/10.1016/j.humov.2005.10.011>
2. World Health Organization. International Classification of Functioning, Disability and Health: ICF: World Health Organization 2001.
3. Magalhaes LC, Cardoso AA, Missiuna C. Activities and participation in children with developmental coordination disorder: a systematic review. *Research in developmental disabilities* 2011;32(4):1309-16. doi: 10.1016/j.ridd.2011.01.029 [published Online First: 2011/02/19]
4. Spittle AJ, Spencer-Smith MM, Eeles AL, et al. Does the Bayley-III Motor Scale at 2 years predict motor outcome at 4 years in very preterm children? *Developmental Medicine And Child Neurology* 2013;55(5):448-52. doi: 10.1111/dmcn.12049
5. Logan SW, Robinson LE, Getchell N. The Comparison of Performances of Preschool Children on Two Motor Assessments. *Perceptual and Motor Skills* 2011;113(3):715-23.
6. Venetsanou F, Kambas A. Environmental factors affecting preschoolers' motor development. *Early Childhood Education Journal* 2010;37(4):319-27.
7. Mokkink LB, Terwee CB, Patrick DL, et al. The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. *Journal of clinical epidemiology* 2010;63(7):737-45. doi: 10.1016/j.jclinepi.2010.02.006 [published Online First: 2010/05/25]
8. Slater LM, Hillier SL, Civetta LR. The clinimetric properties of performance-based gross motor tests used for children with developmental coordination disorder: A systematic review. *PEDIATRIC PHYSICAL THERAPY* 2010;22(2):170.
9. Brown T, Lalor A. The Movement Assessment Battery for Children—Second Edition (MABC-2): A Review and Critique. *Physical & Occupational Therapy in Pediatrics* 2009;29(1):86-103.
10. Terwee CB, Mokkink LB, Knol DL, et al. Rating the methodological quality in systematic reviews of studies on measurement properties: a scoring system for the COSMIN checklist. *Quality of life research : an international journal of quality of life aspects of treatment, care and rehabilitation* 2012;21(4):651-7. doi: 10.1007/s11136-011-9960-1 [published Online First: 2011/07/07]
11. Law M. Outcome measures rating form. Ontario, Canada: CanChild Centre for Disability Research, 2004.
12. Bayley N. Bayley Scales of Infant Development and Toddler Development: Technical Manual: The PsychCorp 2006.
13. Bruininks R, Bruininks B. Bruininks-Oseretsky Test of Motor Proficiency—2nd Edition (BOT-2): Manual. Circle Pines: MN: AGS Publishing 2005.
14. Henderson SE, Sugden DA, Barnett AL. Movement assessment battery for children-2: Movement ABC-2: Examiner's manual: Pearson 2007.

- 1 15. McCarron LT. MAND: McCarron assessment of neuromuscular development, fine and gross motor abilities:
2 McCarron-Dial Systems, Incorporated 1997.
- 3 16. Burns YR. N.S.M.D.A Physiotherapy Assessment for Infants and Young Children Second Edition. Brisbane,
4 Queensland: CopyRight Publishing Company 2014.
- 5 17. Folio M, Fewell R. Peabody Developmental Motor Scales. Examiner's Manual . 2nd Edition. Austin, Texas.:
6 Pro-Ed. 2000.
- 7 18. Ulrich DA. Test of gross motor development-2. *Austin: Prod-Ed* 2000
- 8 19. Rudd J, Butson ML, Barnett L, et al. A holistic measurement model of movement competency in children.
9 *Journal of Sports Sciences* 2016;34(5):477-85.
- 10 20. Valentini NC. Validity and reliability of the TGMD-2 for Brazilian children. *Journal of Motor Behavior*
11 2012;44(4):275-80.
- 12 21. Houwen S, Hartman E, Jonker L, et al. Reliability and validity of the TGMD-2 in primary-school-age children
13 with visual impairments. *Adapt Phys Act Q* 2010;27(2):143-59.
- 14 22. Wong KYA, Yin Cheung S. Confirmatory factor analysis of the Test of Gross Motor Development-2.
15 *Measurement in Physical Education & Exercise Science* 2010;14(3):202-09. doi:
16 10.1080/10913671003726968
- 17 23. Simons J, Daly D, Theodorou F, et al. Validity and reliability of the TGMD-2 in 7-10-year-old Flemish
18 children with intellectual disability. *Adapted physical activity quarterly : APAQ* 2008;25(1):71-82.
- 19 24. Hua J, Gu G, Meng W, et al. Age band 1 of the Movement Assessment Battery for Children-Second Edition:
20 exploring its usefulness in mainland China. *Research in developmental disabilities* 2013;34(2):801-8.
- 21 25. Niemeijer AS, van Waelvelde H, Smits-Engelsman BC. Crossing the North Sea seems to make DCD
22 disappear: cross-validation of Movement Assessment Battery for Children-2 norms. *Hum Mov Sci*
23 2015;39:177-88. doi: 10.1016/j.humov.2014.11.004
- 24 26. Visser L, Ruiters SAJ, Van der Meulen BF, et al. Low verbal assessment with the Bayley-III. *Research in*
25 *developmental disabilities* 2015;36:230-43.
- 26 27. Wuang YP, Su CY. Reliability and responsiveness of the Bruininks-Oseretsky Test of Motor Proficiency-
27 Second Edition in children with intellectual disability. *Research in developmental disabilities* 2009;30(5):847-
28 55.
- 29 28. Wuang YP, Su CY, Huang MH. Psychometric comparisons of three measures for assessing motor functions
30 in preschoolers with intellectual disabilities. *Journal of Intellectual Disability Research* 2012;56(6):567-78.
- 31 29. Ellinoudis T, Evaggelinou C, Kourtessis T, et al. Reliability and validity of age band 1 of the Movement
32 Assessment Battery for Children--second edition. *Research in developmental disabilities* 2011;32(3):1046-51.
33 doi: <http://dx.doi.org/10.1016/j.ridd.2011.01.035>
- 34 30. Smits-Engelsman BCM, Niemeijer AS, van Waelvelde H. Is the Movement Assessment Battery for Children-
35 2nd edition a reliable instrument to measure motor performance in 3 year old children? *Research in*
36 *developmental disabilities* 2011;32(4):1370-77.

- 1 31. Holm I, Tveter AT, Aulie VS, et al. High intra- and inter-rater chance variation of the movement assessment
2 battery for children 2, ageband 2. *Research in developmental disabilities* 2013;34(2):795-800.
- 3 32. Kita Y, Suzuki K, Hirata S, et al. Applicability of the Movement Assessment Battery for Children-Second
4 Edition to Japanese children: A study of the Age Band 2. *Brain & Development* 2016;38(8):706-13.
- 5 33. Griffiths A, Morgan P, Anderson PJ, et al. Predictive value of the Movement Assessment Battery for
6 Children - Second Edition at 4 years, for motor impairment at 8 years in children born preterm. *Dev Med
7 Child Neurol* 2017;59(5):490-96. doi: 10.1111/dmcn.13367 [published Online First: 2017/01/10]
- 8 34. Schulz J, Henderson SE, Sugden DA, et al. Structural validity of the Movement ABC-2 test: factor structure
9 comparisons across three age groups. *Research in developmental disabilities* 2011;32(4):1361-9.
- 10 35. Valentini NC, Ramalho MH, Oliveira MA. Movement assessment battery for children-2: translation,
11 reliability, and validity for Brazilian children. *Research in developmental disabilities* 2014;35(3):733-40. doi:
12 <http://dx.doi.org/10.1016/j.ridd.2013.10.028>
- 13 36. Wuang YP, Su JH, Su CY. Reliability and responsiveness of the Movement Assessment Battery for Children-
14 Second Edition Test in children with developmental coordination disorder. *Developmental Medicine & Child
15 Neurology* 2012;54(2):160-5.
- 16 37. Hands B, Larkin D, Rose E. The psychometric properties of the McCarron Assessment of Neuromuscular
17 Development as a longitudinal measure with Australian youth. *Human Movement Science* 2013;32(3):485-
18 97.
- 19 38. Danks M, Maideen MF, Burns YR, et al. The long-term predictive validity of early motor development in
20 "apparently normal" ELBW survivors. *Early Human Development* 2012;88(8):637-41.
- 21 39. MacDonald J, Burns Y. Performance on the NSMDA During the First and Second Year of Life to Predict
22 Functional Ability at the Age Of 4 in Children with Cerebral Palsy. *Hong Kong Physiotherapy Journal*
23 2005;23(1):40-45. doi: 10.1016/S1013-7025(09)70058-2
- 24 40. Burns YR, Ensbey RM, Norrie MA. The Neuro-sensory motor developmental assessment part 1:
25 development and administration of the test. *Australian Journal of Physiotherapy* 1989;35(3):141-49.
- 26 41. Burns YR, Ensbey RM, Norrie MA. The neuro-sensory motor developmental assessment part II: predictive
27 and concurrent validity. *Australian Journal of Physiotherapy* 1989;35(3):151-57.
- 28 42. Barnett LM, Minto C, Lander N, et al. Interrater reliability assessment using the Test of Gross Motor
29 Development-2. *Journal of Science and Medicine in Sport* 2014;17(6):667-70. doi:
30 <http://dx.doi.org/10.1016/j.jsams.2013.09.013>
- 31 43. Farrokhi A, Zareh Zadeh M, Karimi Alvar L, et al. Reliability and validity of test of gross motor development-
32 2 (Ulrich, 2000) among 3-10 aged children of Tehran City. *Journal of Physical Education and Sports
33 Management* 2014;5(2):18-28. doi: 10.5897/JPESM12.003
- 34 44. Kim Y, Park I, Kang M. Examining rater effects of the TGMD-2 on children with intellectual disability. *Adapt
35 Phys Act Q* 2012;29(4):346-65.
- 36 45. Kim CI, Han DW, Park IH. Reliability and validity of the test of gross motor development-II in Korean
37 preschool children: applying AHP. *Research in developmental disabilities* 2014;35(4):800-7.

- 1 46. Wuang YP, Lin YH, Su CY. Rasch analysis of the Bruininks-Oseretsky Test of Motor Proficiency-Second
2 Edition in intellectual disabilities. *Research in developmental disabilities* 2009;30(6):1132-44.
- 3 47. Spittle AJ, Doyle LW, Boyd RN. A systematic review of the clinimetric properties of neuromotor
4 assessments for preterm infants during the first year of life. *Dev Med Child Neurol* 2008;50(4):254-66. doi:
5 10.1111/j.1469-8749.2008.02025.x [published Online First: 2008/01/15]
- 6 48. McDowell I. Measuring health: a guide to rating scales and questionnaires: Oxford university press 2006.
- 7 49. Cools W, De Martelaer K, Samaey C, et al. Movement skill assessment of typically developing preschool
8 children: a review of seven movement skill assessment tools.(Report). *Journal of Sports Science and Medicine*
9 2009;8(2):154.
- 10 50. Adair B, Said CM, Rodda J, et al. Psychometric properties of functional mobility tools in hereditary spastic
11 paraplegia and other childhood neurological conditions. *Developmental Medicine & Child Neurology*
12 2012;54(7):596-605.

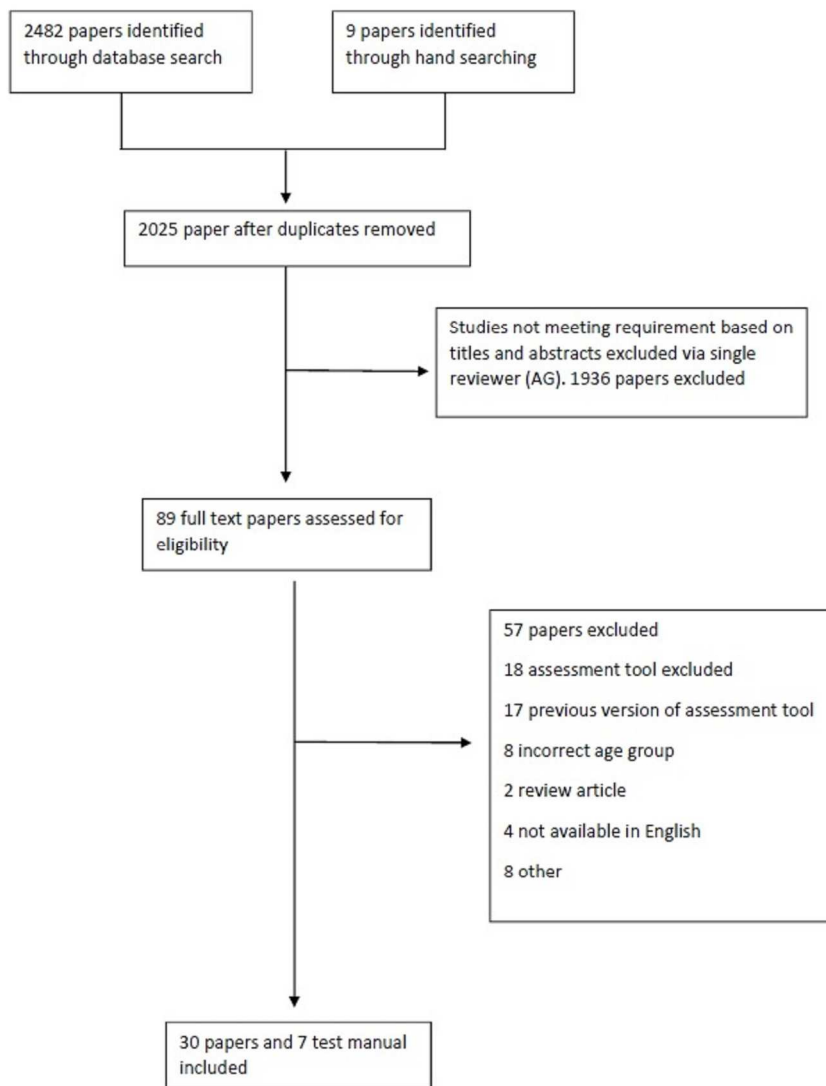


Figure 1. PRISMA flow diagram detailing study selection

195x256mm (300 x 300 DPI)

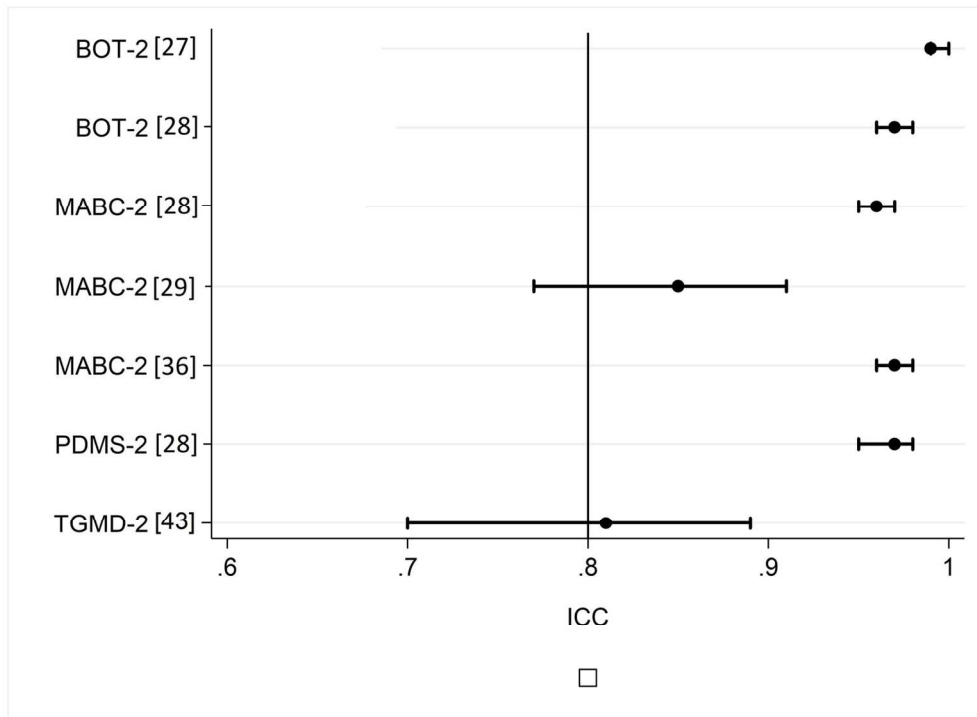


Figure 2. Test re-test reliability of gross motor assessment tools
Figure 2 legend: BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition 13; MABC-2, Movement Assessment Battery for Children 2nd edition 14; PDMS-2, Peabody Developmental Motor Scales 2nd edition 17; TGMD-II, Test of Gross Motor Development 2nd edition 18.

139x102mm (300 x 300 DPI)

Only

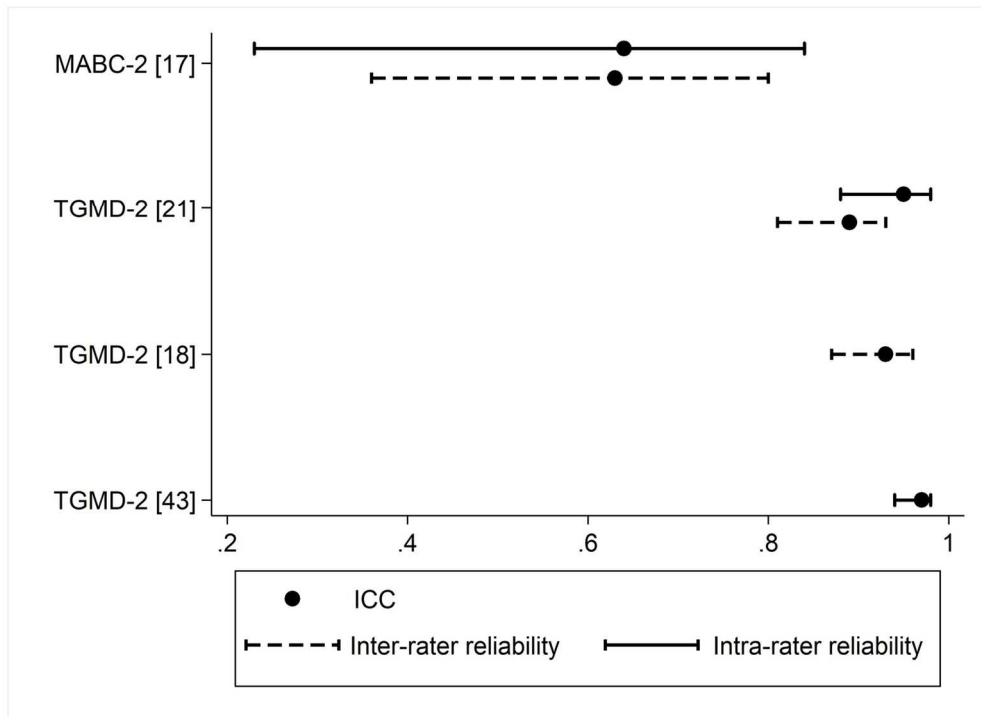


Figure 3. Inter and interrater reliability of gross motor assessment tools
 Figure 3 legend: MABC-2, Movement Assessment Battery for Children 2nd edition 14; TGMD-II, Test of Gross Motor Development 2nd edition 18

139x102mm (300 x 300 DPI)

Supplementary Table 1. OVID Medline database search (1996 to present)

Search No	Search	Yield 5/5/17
1	Child/	811722
2	Child, Preschool/	457484
3	paediatric*.mp.	45528
4	Motor Skills/	12726
5	Motor Activity/	64838
6	gross motor.mp.	3821
7	Psychomotor Disorders/	2609
8	Motor Skills Disorders/	2580
9	Developmental Disabilities/	13484
10	developmental coordination disorder.mp.	845
11	Movement/ph (physiology)	22342
12	Questionnaires/	336296
13	"Outcome Assessment (Health Care)"/	57491
14	scale*.mp.	608566
15	instrument*.mp.	197131
16	outcome*.mp.	1813266
17	measure*.mp.	2255187
18	evaluat*.mp.	2552240
19	assess*.mp.	2273012
20	"Task Performance and Analysis"/	22017 (or 5969)
21	Reproducibility of Results"/	319899
22	1 or 2 or 3	936097
23	4 or 5 or 6 or 7 or 8 or 9 or 10 or 11	116200
24	12 or 13 or 14 or 15 or 16 or 17 or 18 or 19 or 20	6579985
25	21 and 22 and 23 and 24	1152

Supplementary table 2: CINAHL plus database search

Search Number	Search	Yield
		2/7/17
S1	(MH "Child")	338732
S2	(MH "Child, Preschool")	156847
S3	"paediatric"	14351
S4	(MH "Motor Skills")	7420
S5	(MH "Motor Activity")	9664
S6	(MH "Psychomotor Performance")	9457
S7	(MH "Motor Skills Disorders")	1515
S8	(MH "Developmental Disabilities")	7114
S9	(MH "Child Development Disorders")	1708
S10	"gross motor"	2234
S11	(MH "Clinical Assessment Tools")	110291
S12	(MH "Outcome Assessment")	29335
S13	(MH Physical Therapy Assessment")	2027
S14	"scale"	290924
S15	instrument*	113984
S16	outcome*	602132
S17	measure*	534465
S18	evaluat*	760447
S19	assess*	710533
S20	(MH Reliability and Validity")	11043
S21	S1 OR S2 OR S3	387809
S22	S4 OR S5 OR S6 OR S7 OR S8 OR S9 OR S10	35562
S23	S11 OR S12 OR S13 OR S14 OR S15 OR S16 OR S17 OR S18 OR S19	1763829
S24	S20 AND S21 AND S22 AND S23	144

Supplementary table 3: EMBASE database search (1974-present)

Search No	Search	Search Yield 5/5/17
1	Child/	1346641
2	Preschool child/	457981
3	paediatric*.mp.	86036
4	Motor performance/	57571
5	Motor Activity/	39751
6	psychomotor performance	19515
7	Motor development/	4906
8	Motor dysfunction/	53155
9	Developmental disorder/	30473
10	Gross motor.mp.	6840
11	Outcome Assessment/	358121
12	Outcome measure.mp.	60507
13	Questionnaire/	513199
14	Task performance/	125167
15	Functional assessment/	55415
16	Clinical assessment tool/	19865
17	evaluat*.mp.	3874341
18	instrument*.mp.	515930
19	outcome*.mp.	2425627
20	Assess*.mp.	3815907
21	Scale*.mp.	903216
22	Measure*.mp.	3444366
23	Measurement accuracy/	18209
24	Measurement repeatability/	2849
25	Reproducibility/	173988
26	Validity/	40192
27	Reliability/	114002
28	1 or 2 or 3	1535605
29	4 or 5 or 6 or 7 or 8 or 9 or 10	195237
30	11 or 12 or 13 or 14 or 15 or 16 or 17 or 18 or 19 or 20 or 21 or 22	10121631
31	23 or 24 or 25 or 26 or 27	324779
32	28 and 29 and 30 and 31	1105

Supplementary table 4. Allied and Complementary Medicine Database (AMED) database search: (1985-present)

Search Number	Search	Yield 2/7/17
1	Child/	15192
2	Child preschool/	1223
3	Adolescent/	3979
4	paediatric*.mp.	812
5	1 or 2 or 3 or 4	18429
6	Motor skills/	1220
7	Motor activity/	1468
8	Gross motor*.mp.	599
9	Psychomotor disorders/	1067
10	Developmental disabilities/ or motor skills disorders/	947
11	Developmental coordination disorder*.mp.	219
12	DCD.mp.	113
13	6 or 7 or 8 or 9 or 10 or 11 or 12	4982
14	5 and 13	1510
15	Clinical assessment scales/	4318
16	Questionnaires/	4123
17	Disability evaluation/	7023
18	Outcome measure*.mp/	9845
19	Outcome*.mp.	38379
20	Assess*.mp.	43680
21	Scale*.mp.	17562
22	Evaluat*.mp.	40621
23	15 or 16 or 17 or 18 or 19 or 20 or 21 or 22	93570
24	14 and 23	865
25	Measurement/	1629
26	Reproducibility of results.mp/	2241
27	"Consistency and reliability"/	1898
28	Statistics/	1075
29	Specificity.mp.	1241
30	Sensitivity.mp.	2860
31	"Predictive value of tests"/	839
32	25 or 26 or 27 or 28 or 29 or 30 or 31	10256
33	24 and 32	81

Supplementary table 5: Definition of terms

	Measurement Property	Definition	Example/explanation
Validity	Content	The degree to which an assessment tool's content measures the construct that it intends to measure ⁷	Concerned with the relevance and comprehensiveness of the items included in the assessment tool
	Construct	Measures the degree to which the scores obtained from the test are an adequate reflection of the construct to be measured ⁷	Examples include structural validity (whether scores reflect the dimensionality of the construct), hypothesis testing (item construct validity) and cross-cultural validity (whether translated or culturally adapted assessments adequately reflect the original version) ⁷
	Criterion	Assesses whether or not the test scores reflect a 'gold standard' assessment ⁷	As there is no gold standard of assessment for gross motor function in children this is often assessed with correlations of scores obtained from two or three other frequently used tools.
Reliability	Reliability	Refers to the consistency of a test score regardless of the time between assessments (test-retest) or the person administering (intra and inter-rater) ⁵⁰	Usually measured with intraclass correlation coefficient (ICC), but can be measured using Cohen's kappa coefficient. Percentage agreement and Pearson's correlation coefficient do not incorporate error into the calculations and as such is not a true measure of agreement ⁵⁰ . Scores > 0.80 are considered excellent, 0.60-0.79 adequate and <0.59 poor ¹¹
	Internal consistency	The degree of interrelatedness of an assessment tool's items ⁷	Usually measured using Cronbach's alpha (α) ⁷ . scores > 0.70 demonstrates high relationship, 0.5 to 0.69 a moderate relationship, 0.26 to 0.49 a low relationship and < 0.26 little relationship ⁵⁰ .
	Measurement Error	Refers to the error obtained between measurements that cannot be attributed to the patients true change ⁷	May be systematic or random error ⁷
Responsiveness	Responsiveness	An assessment tool's ability to detect change over time in the construct it purports to measure ⁷	This is central to a tools capacity to be used as an outcome measure.

Supplementary table 6: Excluded Assessment Tools

Reason	Assessments
Manual not available in English	Maastricht's Motor Test (MMT) The Motor-Proficiency-Test for children between 4 and 6 years of age (MOT 4-6) Zuk Assessment Körperkoordinationstest für Kinder (KTK)
Cannot extract meaningful gross motor score	Early Intervention Developmental Profile (EIDP) Neurological Developmental Exam Preschooler Gross Motor Quality Scale (PGMQ) The Malawi Developmental Assessment Tool (MDAT) Dutch table tennis motor skills assessment
Screening Tool	Brief Assessment of Motor Function (BAMF) The Motor Performance Checklist Motor skill checklist (MSC)
Diagnosis specific/requires a diagnosis	Assessment Battery for the Atypical Handicapped Child (VAB) Video-based documentation and rating system of the motor behaviour of handicapped children
Only assesses one motor domain (e.g. gait)	Standardized Walking Obstacle Course (SWOC) Timed floor to stand test
Manual not published/commercially available	Rapid Neurodevelopmental Assessment (RNDA) Tufts Assessment of Motor Performance (TAMP) Zurich Neuromotor Assessment (ZNA)

Supplementary table 7: Scoring and administration of assessment tools

Assessment Tool	Scoring	Interpretation of scores	Other
Bayley-III ¹²	Motor score - gross (varying items) and fine motor (varying items) subscales. Binary score with reverse/discontinue rules	Raw scores Composite scores Centile ranks Age equivalents Growth scores	Lends itself to multidisciplinary team testing.
BOT-2 ¹³	Fine manual (15 items) manual coordination (12 items) body coordination (16 items) strength and agility (10 items) subscales. Scoring differs for subtests	Raw scores Age adjusted standard scores Composite scores Centile ranks Age equivalents Descriptive categories. Complex conversions	Administration Easel includes instructions, diagrams and photos of test procedure
MABC-2 ¹⁴	Manual dexterity (3 items), aiming & catching (2 items) and balance (3 items) subscales.	Raw scores component scores centile ranks total test score traffic light system. Simple conversion	Also Available: MABC-2 Checklist (screening tool) and intervention manual
MAND ¹⁵	Fine motor (5 items) Gross motor (5 items)	Raw scores Scaled scores converted to an NDI. Factor scores. Complex conversions	Case studies included in manual for hyperactivity, encephalitis, mild head trauma, CP and muscular dystrophy
NSMDA ¹⁶	Functional grade given for each subscale, which is combined to create an overall score.	Indicates: normal range, minimal dysfunction, mild problems, moderate, severe or profound disability	Sections for comment on strengths, behavioural state during testing, musculoskeletal system and recommendations.
PDMS-2 ¹⁷	GM: Stationary (30 items), locomotion (89 items), object manipulation (24 item). FM: grasping(26 items) , visual-motor integration (72 items)	Raw scores, Age equivalent, centile rank. Standard scores (subtests) Composite quotient. Complex conversions.	Motor activities program (intervention ideas)
TGMD-2 ¹⁸	Locomotion (6 items) and Object Control (6 items). Separate male/female norms for object control subset	Raw scores, standard scores, percentile rank, age equivalent, Gross Motor Quotient. Simple conversion.	Simple to administer

Bayley-III, Bayley Scale of Infant and Toddler Development 3rd edition ¹²; BOT-2, Bruininks-Oseretsky Test of Motor Proficiency 2nd edition ¹³; MABC-2, Movement Assessment Battery for Children 2nd edition ¹⁴; MAND, McCarron Assessment of Neuromuscular Development ¹⁵; NSMDA, Neurological Sensory Motor Developmental Assessment ¹⁶; PDMS-2, Peabody Developmental Motor Scales 2nd edition ¹⁷; TGMD-II, Test of Gross Motor Development 2nd edition ¹⁸; GM, Gross Motor; FM, Fine Motor; NDI, Neurodevelopmental Index

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only



PRISMA 2009 Checklist

Section/topic	#	Checklist item	Reported on page #
TITLE			
Title	1	Identify the report as a systematic review, meta-analysis, or both.	1 Title page
ABSTRACT			
Structured summary	2	Provide a structured summary including, as applicable: background; objectives; data sources; study eligibility criteria, participants, and interventions; study appraisal and synthesis methods; results; limitations; conclusions and implications of key findings; systematic review registration number.	2
INTRODUCTION			
Rationale	3	Describe the rationale for the review in the context of what is already known.	4-5
Objectives	4	Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS).	5
METHODS			
Protocol and registration	5	Indicate if a review protocol exists, if and where it can be accessed (e.g., Web address), and, if available, provide registration information including registration number.	-
Eligibility criteria	6	Specify study characteristics (e.g., PICOS, length of follow-up) and report characteristics (e.g., years considered, language, publication status) used as criteria for eligibility, giving rationale.	5-6
Information sources	7	Describe all information sources (e.g., databases with dates of coverage, contact with study authors to identify additional studies) in the search and date last searched.	5
Search	8	Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated.	5 and suppl. tables 1, 2, 3, 4
Study selection	9	State the process for selecting studies (i.e., screening, eligibility, included in systematic review, and, if applicable, included in the meta-analysis).	5-6
Data collection process	10	Describe method of data extraction from reports (e.g., piloted forms, independently, in duplicate) and any processes for obtaining and confirming data from investigators.	6-7
Data items	11	List and define all variables for which data were sought (e.g., PICOS, funding sources) and any assumptions and simplifications made.	6-7
Risk of bias in individual studies	12	Describe methods used for assessing risk of bias of individual studies (including specification of whether this was done at the study or outcome level), and how this information is to be used in any data synthesis.	6
Summary measures	13	State the principal summary measures (e.g., risk ratio, difference in means).	7
Synthesis of results	14	Describe the methods of handling data and combining results of studies, if done, including measures of consistency (e.g., I^2) for each meta-analysis.	-



PRISMA 2009 Checklist

Page 1 of 2

Section/topic	#	Checklist item	Reported on page #
Risk of bias across studies	15	Specify any assessment of risk of bias that may affect the cumulative evidence (e.g., publication bias, selective reporting within studies).	24, 25
Additional analyses	16	Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified.	-
RESULTS			
Study selection	17	Give numbers of studies screened, assessed for eligibility, and included in the review, with reasons for exclusions at each stage, ideally with a flow diagram.	Figure 1 + page 7
Study characteristics	18	For each study, present characteristics for which data were extracted (e.g., study size, PICOS, follow-up period) and provide the citations.	Table 1 – page 9 + Suppl table 7
Risk of bias within studies	19	Present data on risk of bias of each study and, if available, any outcome level assessment (see item 12).	11 + Table 3 (page 14-15)
Results of individual studies	20	For all outcomes considered (benefits or harms), present, for each study: (a) simple summary data for each intervention group (b) effect estimates and confidence intervals, ideally with a forest plot.	7-8, 11-12, 17 + Table 2 page 13, Table 4 page 18, Table 5 page 19, Table 6 page 20 + Figures 2 & 3
Synthesis of results	21	Present results of each meta-analysis done, including confidence intervals and measures of consistency.	-
Risk of bias across studies	22	Present results of any assessment of risk of bias across studies (see Item 15).	-
Additional analysis	23	Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]).	-
DISCUSSION			
Summary of evidence	24	Summarize the main findings including the strength of evidence for each main outcome; consider their relevance to key groups (e.g., healthcare providers, users, and policy makers).	22-25



PRISMA 2009 Checklist

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47

Limitations	25	Discuss limitations at study and outcome level (e.g., risk of bias), and at review-level (e.g., incomplete retrieval of identified research, reporting bias).	24
Conclusions	26	Provide a general interpretation of the results in the context of other evidence, and implications for future research.	26
FUNDING			
Funding	27	Describe sources of funding for the systematic review and other support (e.g., supply of data); role of funders for the systematic review.	1

From: Moher D, Liberati A, Tetzlaff J, Altman DG, The PRISMA Group (2009). Preferred Reporting Items for Systematic Reviews and Meta-Analyses: The PRISMA Statement. PLoS Med 6(7): e1000097. doi:10.1371/journal.pmed1000097

For more information, visit: www.prisma-statement.org.