

International Journal of Health Professions

A systematic review of assessments for procedural skills in physiotherapy education.

--Manuscript Draft--

Manuscript Number:	IJHP-D-16-00023R1
Article Type:	Original Study
Full Title:	A systematic review of assessments for procedural skills in physiotherapy education.
Short Title:	A systematic review of assessments for procedural skills in physiotherapy education.
Secondary Full Title:	Assessment von prozeduralen Fähigkeiten in der physiotherapeutischen Ausbildung: Ein systematischer Review
Secondary Short Title:	Assessment von prozeduralen Fähigkeiten in der physiotherapeutischen Ausbildung: Ein systematischer Review
Keywords:	procedural skills, practical skills, systematic review, clinical assessment
Secondary Keywords:	Procedurale Fähigkeiten, praktische Fähigkeiten, systematischer Review, klinisches Assessment
Abstract:	<p>Introduction: Learning of procedural skills is important in the education of physiotherapists. It is the aim of physiotherapy degree programmes that graduates are able to practice selected procedures safely and efficiently. Procedural competency is threatened by an increasing and diverse amount of procedures that are incorporated in university curricula. As a consequence, less time is available for the learning of each specific procedure. Incorrectly performed procedures in physiotherapy might be ineffective and may result in injuries to patients and physiotherapists. The aim of this review was to synthesise relevant literature systematically to appraise current knowledge relating to assessments for procedural skills in physiotherapy education.</p> <p>Method: A systematic search strategy was developed to screen five relevant databases (CINAHL, Cochrane Central, SportDISCUS, ERIC and MEDLINE) for eligible studies. The included assessments were evaluated for evidence of their reliability and validity.</p> <p>Results: The search of electronic databases identified 560 potential records. Seven studies were included into this systematic review. The studies reported eight assessments of procedural skills. Six of the assessments were designed for a specific procedure and two assessments were considered for the evaluation of more than one procedure. Evidence to support the measurement properties of the assessment was not available for all categories.</p> <p>Discussion: It was not possible to recommend a single assessment of procedural skills in physiotherapy education following this systematic review. There is a need for further development of new assessments to allow valid and reliable assessments of the broad spectrum of physiotherapeutic practice.</p>
Secondary Abstract:	<p>Einleitung: Das Erlernen von prozeduralen Fähigkeiten ist ein wichtiges Element in der Ausbildung von Physiotherapeuten. Es ist das Ziel von physiotherapeutischen Studiengängen, dass Graduierte in der Lage sind, ausgewählte Prozeduren sicher und effektiv auszuführen. Die prozedurale Kompetenz ist bedroht von wechselnden und einer stetig anwachsenden Anzahl von Prozeduren, die in die Curricula der Studiengänge eingebaut werden. Als Konsequenz ist weniger Zeit vorhanden, um die einzelnen Prozeduren zu erlernen. Falsch durchgeführte Prozeduren können zu Verletzungen von Patienten und Physiotherapeuten führen.</p> <p>Zielsetzung der Arbeit war es, relevante Literatur systematisch zu erfassen, um eine Übersicht von Assessments von prozeduralen Fähigkeiten in der physiotherapeutischen Ausbildung zu erstellen.</p> <p>Methode: Eine systematische Suchstrategie wurde entwickelt, um fünf Datenbanken (CINAHL, Cochrane Central, SportDISCUS, ERIC and MEDLINE) nach relevanten Studien zu durchsuchen. Die eingeschlossenen Assessments wurden im Bezug auf ihre Reliabilität und Validität bewertet.</p> <p>Ergebnisse: Die Suche in den elektronischen Datenbanken ergab 560 Treffer. Sieben Studien wurden in diese systematische Übersichtsarbeit eingeschlossen. Die Studien berichteten über acht Assessments für prozedurale Fähigkeiten. Sechs Assessments sind für eine spezifische Prozedur entwickelt worden und zwei Assessments können</p>

	für unterschiedliche Prozeduren benutzt werden. Evidenz für die Messeigenschaften der eingeschlossenen Messinstrumente war nicht für alle Kategorien verfügbar. Diskussion: Es ist nicht möglich, ein bestimmtes Messinstrument zur Bewertung von prozeduralen Fähigkeiten zu empfehlen. Es gibt einen Bedarf an Messinstrumenten, die reliabel und valide sind, um das breite Spektrum von prozeduralen Fähigkeiten zu bewerten.
Corresponding Author:	Martin Sattelmayer HES-SO Valais Wallis Leukerbad, SWITZERLAND
First Author:	Martin Sattelmayer
Order of Authors:	Martin Sattelmayer
	Roger Hilfiker
	Gillian Baer
Manuscript Region of Origin:	SWITZERLAND

1 A systematic review of assessments for
2 procedural skills in physiotherapy
3 education

4

5 Keywords: procedural skills, practical skills, systematic review, clinical assessment

6

7 **1 Abstract**

8 Introduction: Learning of procedural skills is important in the education of physiotherapists.

9 It is the aim of physiotherapy degree programmes that graduates are able to practice
10 selected procedures safely and efficiently. Procedural competency is threatened by an
11 increasing and diverse amount of procedures that are incorporated in university curricula. As
12 a consequence, less time is available for the learning of each specific procedure. Incorrectly
13 performed procedures in physiotherapy might be ineffective and may result in injuries to
14 patients and physiotherapists. The aim of this review was to synthesise relevant literature
15 systematically to appraise current knowledge relating to assessments for procedural skills in
16 physiotherapy education.

17 Method: A systematic search strategy was developed to screen five relevant databases
18 (CINAHL, Cochrane Central, SportDISCUS, ERIC and MEDLINE) for eligible studies. The
19 included assessments were evaluated for evidence of their reliability and validity.

20 Results: The search of electronic databases identified 560 potential records. Seven studies
21 were included into this systematic review. The studies reported eight assessments of
22 procedural skills. Six of the assessments were designed for a specific procedure and two
23 assessments were considered for the evaluation of more than one procedure. Evidence to
24 support the measurement properties of the assessment was not available for all categories.

25 Discussion: It was not possible to recommend a single assessment of procedural skills in
26 physiotherapy education following this systematic review. There is a need for further
27 development of new assessments to allow valid and reliable assessments of the broad
28 spectrum of physiotherapeutic practice

29

30 2 Introduction

31 It is the aim of physiotherapy degree programmes that graduates are able to execute
32 selected procedures safely and efficiently. Considerable resources are allocated to enable
33 graduates to achieve a high level of procedural competency. Within this review procedural
34 skills were classified after Kent's definition as: "a skill involving a series of discrete responses
35 each of which must be performed at the appropriate time in the appropriate sequence"
36 (Kent, 2007, p. 437).

37 Recent literature highlights that there is no consensus with regard to definitions and
38 classifications of procedural skills. Michels, Evans, and Blok (2012) identified that procedural
39 skills are not exactly defined in the field of health professions education. Frequently, they
40 are categorised under the umbrella term "clinical skills". However, there is a lack of
41 standardisation. Simpson et al. (2002), separated practical procedures from communication
42 skills, clinical skills, and other skills in the Scottish doctor learning outcomes. In contrast, the
43 General Medical Council in the UK does not separate between procedural skills and clinical
44 skills (2004), for example safety measures are categorised as essential procedural skills in
45 their classification. Lastly, the Royal Australian College of General Practitioners (2011)
46 defined procedural skills as: "A procedure is a manual intervention that aims to produce a
47 specific outcome during the course of patient care" (The Royal Australian College of General
48 Practitioners, 2011, p . 515).

49 To avoid ambiguity in this review, procedural skills were characterised with the following
50 features: a) they involve the execution of a procedural task (e.g. a manual or a practical
51 task), b) involvement of technical equipment may be possible but this is not a prerequisite of
52 procedural skills, c) the character of a procedure can be diagnostic, evaluative or
53 interventional and d) procedures can range from simple tasks with few parts to complex
54 sequences involving multiple activities.

55 As procedures in physiotherapy are highly interactive between patients and therapists, more
56 information than execution of procedures may be needed to evaluate procedural skills. For
57 example, communication providing basic information about the procedures between
58 physiotherapist and patient is frequently necessary. Consequently, therapists should be

59 educated to allow them to adapt procedures to a variety of circumstances such as
60 environmental requirements or individual patient needs.

61 Physiotherapy is a dynamic profession with evolution of new physiotherapeutic roles and
62 skills in many health systems (Higgs, Hunt, Higgs, & Neubauer, 1999) thus requiring the
63 incorporation of new tasks and skills into physiotherapy degree curricula. However, this may
64 result in an increased amount of procedures that are incorporated in university curricula. As
65 a consequence, less time is available for the learning of specific procedures.

66 Incorrectly performed procedures in physiotherapy might be ineffective and may result in
67 injuries to physiotherapists or to patients. For example, Nyland and Grimmer (2003)
68 reported that low back pain is frequently experienced by undergraduate physiotherapy
69 students and, Glista and co-workers (2014) reported that the students' posture deteriorated
70 during the course of education. In some situations, physiotherapists are required to perform
71 professional procedures in difficult environments with poor working postures which are
72 potential harmful for the musculoskeletal system (Jackson & Liles, 1994). Therefore, training
73 of procedures should be designed to enable learners to perform procedures without
74 endangering their own personal safety and to understand how to adapt procedures
75 appropriately.

76 Procedures performed by physiotherapist can also be associated with adverse events for
77 patients. For example, Gorrell, Engel, Brown, and Lystad (2016) reported that mild adverse
78 events occurred in 61 RCTs and major adverse events were seen in 2 RCTs evaluating spinal
79 manipulative therapy. Therefore, following the initial teaching of procedural skills,
80 physiotherapy educators need valid and reliable assessment tools to evaluate whether
81 procedural competency of students is sufficient for practice.

82 Assessment of procedural skills has been extensively researched in surgical education
83 (Jelovsek, Kow, & Diwadkar, 2013). Some assessments exists, which can be used for
84 procedures in nursing education (Morris, Gallagher, & Ridgway, 2012). While teaching of
85 procedural skills is a core part of undergraduate physiotherapy education, no review could
86 be identified of assessment tools for procedural skills in physiotherapy education.

87 One important consideration in the evaluation of procedural skills in physiotherapy is
88 whether an assessment framework exists. Miller (1990) argued that no single assessment
89 would be sufficient to allow the judgement of such complex skills. He presented a four level

90 framework for assessments in health professions education. The base of this framework is
91 knowledge (the student “knows”), which can be tested with standardised objective test
92 methods (e.g. multiple choice tests). The second level (competence) provides evidence that
93 students know how to use their knowledge (e.g. vignette assessments). The third level
94 evaluates the performance of students (e.g. students have to show how they perform a
95 specific procedure). Lastly, the question remains whether the learned skills are
96 independently selected and used appropriately in clinical practice. Examples to evaluate the
97 “action level” are work place based assessments or portfolios (Chandratilake, Davis, &
98 Ponnampereuma, 2010).

99 The aim of this review was to identify, examine and synthesise relevant literature to produce
100 a systematic review of assessments for procedural skills in physiotherapy education.
101 Specifically, the objective of this review was to identify existing assessments of procedural
102 skills in physiotherapy education and to evaluate them with regard to their measurement
103 properties.

104 **3 Methods**

105 A systematic review was undertaken to address the identified objectives. To increase clarity
106 of reporting, the PRISMA guideline was followed (Liberati et al., 2009).

107 **3.1 Criteria for in and exclusion**

108 Inclusion and exclusion criteria are presented in Table 1.

109 [Table 1. In-and exclusion criteria](#)

110

111 **3.2 Search methods**

112 Five electronic databases were systematically searched for potential eligible studies. These
113 databases were: Cumulative Index to Nursing and Allied Health Literature (CINAHL),
114 Cochrane Central Register of Controlled Trials (CENTRAL), SPORTDiscus, Educational
115 Resource Information Center (ERIC) and Medline via Pubmed. In addition, the references of
116 all included full text articles were checked for relevant studies. The search string is presented
117 in Table 2. Findings of the three categories Population, Assessment and Outcome were
118 combined with the Boolean operator AND.

119 [Table 2. Search strategy](#)

120

121 All retrieved records were imported into an electronic database and duplicates were
122 removed. In a next step, titles and abstracts of the records were screened with regard to the
123 pre-specified inclusion and exclusion criteria. Lastly, the full texts of the remaining studies
124 were read and studies were included in the systematic review if they met all criteria.

125 **3.3 Data collection and management**

126 Data were extracted in relation to the following information:

- 127 • Study details (country, setting and sample)
- 128 • Assessment characteristics (name of the assessment, assessment items, assessment
129 aim, assessment duration, assessment criteria, assessors, patients and target
130 procedure)

- 131 • Measurement properties (internal consistency, reliability, measurement error,
132 content validity and construct validity)
- 133 • Methodological quality of assessments (the Standards for Evaluating the Quality of
134 Assessment Methods in Medical Education (Swing, Clyman, Holmboe, & Williams,
135 2009)

136 3.4 Analysis

137 Evidence of reliability and validity of the included assessments was evaluated. Within
138 reliability the internal consistency, the inter- and intrarater reliability and the measurement
139 error were appraised. Validity was appraised with regard to content validity, criterion
140 validity and construct validity. Despite some discussion about agreed definitions regarding
141 measurement properties, the consensus definitions proposed by Mokkink et al. (2010) were
142 used to ensure consistency in how findings were interpreted.

143 3.5 Assessment of methodological quality of assessments

144 All included assessments were evaluated with the Standards for Evaluating the Quality of
145 Assessment Methods (SEQAM) (Swing et al., 2009). The SEQAM is an assessment tool for
146 educational assessments specifically designed for health professions education. The SEQAM
147 critically evaluates 6 dimensions: reliability (e.g. reliability indicators are available for all used
148 scores), validity (e.g. selection of content is justified), ease of use (e.g. the tool is easily
149 carried out in daily practice), resources required (e.g. training requirements for assessors do
150 not exceed one hour), ease of interpretation (e.g. individual scores are interpretable) and
151 educational impact (e.g. provides useful results). For each dimension studies could be rated
152 as evidence level A, B, C or not rated. For an evidence level of A all standards of one
153 dimension had to be met. Studies were rated as evidence level B when one standard was not
154 met. When two standards in one dimension were not met an evidence level of C was
155 specified. Lastly, when three or more standards were not met an evidence level of not rated
156 (NR) was given. The scoring rules of the SEQAM were adapted from Swing et al. (2009).

157 4 Results

158 The results of this review are presented in three sections. First the results of the search are
159 presented, then findings of the measurement properties of the included assessments are
160 provided. Finally, the methodological quality of the included assessments is considered.

161 4.1 Results of the search

162 The search of electronic databases identified 560 potential records. Additionally, 10 articles
163 were identified by reference checking. It was possible to delete 6 duplicates. Therefore, titles
164 and abstracts of 564 records were screened. The majority of 454 records were excluded
165 because they did not report an appropriate assessment (n= 387). Fifty records did not report
166 an appropriate outcome and 17 records did not meet the inclusion criteria with regard to
167 the population.

168 110 full-text articles were then read. It was possible to exclude 103 full-text articles. Most
169 studies (n = 93) were excluded because they were related to a different discipline in
170 medicine (e.g. surgery). Two studies had insufficient data to include them into the
171 systematic review. They evaluated multiple different patient encounters and therefore it
172 was not possible to extract data for a single assessment method. Eight studies were not
173 included because they were reviews of primary studies. Finally, seven studies were included
174 into this systematic review. The studies reported six procedure specific measurement
175 instruments (PSMI) and two procedure unspecific measurement instruments (PUMI) (Figure
176 1).

177

178 [Figure 1. Study flow](#)

179

180 4.1.1 Included assessments

181 The included assessments were classified as either procedure specific measurement
182 instruments (i.e. assessments designed for one specific procedure) or procedure unspecific
183 measurement instruments (i.e. generic assessments, which can be used for more than one
184 specific procedure).

185 4.1.1.1 Procedure specific measurement instruments

186 The six PSMIs included in this review are briefly presented below. A detailed critical
187 overview is presented in Table 3. The Assessment of Musculoskeletal Physical Examination
188 Skills Checklist (AMPE) was published by Beran et al. (2012). The AMPE is a 12-15 item
189 checklist and evaluates the ability of health professionals to perform a physical examination
190 of four different clinical scenarios. The scenarios involve an upper extremity, a trauma, a
191 spine and a lower extremity case. The AMPE requires in addition to an assessor a trained
192 standardised patient for each of the four scenarios. The authors designed checklists of
193 important procedures, which students should perform when they encounter a specific
194 simulated patient, such as joint palpation or strength testing.

195 Herbers, Wessel, El-Bayoumi, Hassan, and St Onge (2003) created the 29-item Pelvic
196 Examination Skills Checklist (PES-C) and the 5-point Pelvic Examination Skill Rating Scale
197 (PES-R). Most of the 29 items on the PES-C are related to the physical performance of a
198 pelvic examination, although some of the items relate to communication skills (e.g. item 21:
199 Tells patient to state if pain too great). The PES-R is a five-point global rating scale that
200 enables the evaluator to rate the overall performance of the pelvic examination. Both
201 assessments were validated with gynaecologic teaching associates who fulfilled a dual role
202 as subjects for the pelvic examination and evaluators of the learner's performance within
203 the study of Herbers and colleagues.

204 The Physical Examination Skills Checklist (PhyES) was published by Ladyshevsky, Baker,
205 Jones, and Nelson (2000) and aims to evaluate a musculoskeletal physical examination of a
206 patient with a rotator cuff problem. The PhyES is scored on a three-point system and uses
207 carefully coached persons to portray specific patients. Performance was scored using a
208 checklist which included important features of the physical examination (e.g. evaluation of
209 shoulder girdle stability).

210 Swift and colleagues (2013) designed the mOSCE-Station 3 checklist (mO-S3). The mO-S3
211 evaluates the ability of physiotherapy students to perform two specific shoulder assessment
212 tests. Learners have to choose two tests to confirm their hypothesis with regard to a
213 scenario with a patient suffering from shoulder pain. The mO-S3 consists of five
214 dichotomous items and one ordinal item. In order to administer the mO-S3 standardised
215 patients and specialised clinical instructors are necessary. The following tasks were

216 evaluated in the OSCE: i) think station, ii) explanation of the primary hypothesis to a patient,
217 iii) performing two specific tests to confirm the hypothesis, iv) performing the best day 1
218 hands-on intervention, v) reassessment, vi) performing the best day 1 exercise intervention
219 and vii) performing a specific technique and explanation of the selected technique.

220 The 138 item checklist head- to-toe physical examination checklist (HTTPE) (Yudkowsky et
221 al., 2004) evaluates the ability of an “assessor” to perform a complete physical screening
222 examination of the whole body and all 138 items are scored on a trichotomous scoring
223 system. To administer the HTTPE, trained standardised patient instructors are required. The
224 patient instructors serve as patients and mark the “assessors” performance.

225 4.1.1.2 Procedure unspecific measurement instruments

226 The Osteopathic Manipulative Treatment assessment tool (OMT) (Boulet, Gimpel, Dowling,
227 & Finley, 2004) aims to measure the ability to perform a manipulative treatment and
228 consists of 15 items scored on a trichotomous scale. It can be used for different manipulative
229 treatment techniques and for different body regions and therefore is procedure unspecific.
230 For example, Boulet et al. (2004) used the OMT to evaluate various procedures related to
231 the treatment of low back pain, frozen shoulder or asthma. Standardised patients are a
232 prerequisite to use the OMT as an assessment tool.

233 The Global Procedural Skills Evaluation Form (GPSE) was originally presented in the field of
234 family medicine (Nothnagle, Reis, Goldman, & Diemers, 2010). However, its generalised
235 design as a rating scale for procedural skills affords its utility for assessment of procedural
236 skills in physiotherapy as well. The GPSE provides feedback based on direct observation of a
237 procedure. The scoring system is based on a 4-point scale and quantifies the amount of
238 guidance that was needed to perform a procedure. No standardised patients are required
239 when the GPSE is applied. Furthermore, student’s self-assessment is included in the GPSE
240 score.

241

242 **Table 3. Characteristics of included studies and assessments**

243

244 4.2 Findings

245 Within this section evidence of the measurement properties of the included assessments are
246 presented. The consensus definitions proposed by Mokkink et al. (2010) were used to
247 appraise the measurement properties.

248 4.2.1 Reliability

249 Reliability of the assessments was appraised with regard to their internal consistency,
250 interrater reliability, intrarater reliability and measurement error.

251 Two studies were included that reported the internal consistency of two different
252 assessments. Swift et al. (2013) reported an internal consistency between $\alpha = 0.31$ (video
253 examiner) and $\alpha = 0.55$ (onsite examiner) for the mO-S3. They calculated the internal
254 consistency of a 6 station OSCE. The statistical method used to calculate the internal
255 consistency was Cronbach's alpha. Boulet et al. (2004) reported an internal consistency for
256 the OMT between 0.83 (Case 1: low back pain) and 0.97 (Case 3: asthma). All internal
257 consistency estimates are presented in Figure 2.

258

259 **Figure 2. Internal consistency estimates.**

260 **Nb. The statistical method from Boulet et al. (2004) was not available.**

261

262 Six studies were included that reported interrater reliability of six assessments. Beran et al.
263 (2012) evaluated four different procedures using the AMPE. Interrater reliability ranged
264 between 0.27 (95%CI: 0 to 0.56) for the physical examination of trauma patients to 0.77
265 (95% CI: 0.46 to 0.9) for a physical examination of the knee.

266 Herbers et al. (2003) investigated the interrater reliability of students performing a specific
267 pelvic examination with no deviations from the protocol allowed and reported kappa
268 coefficient of $\kappa = 0.54$ for the PES-C (pelvic examination).

269 Ladyshewsky et al. (2000) investigated the interrater reliability for the assessment of a
270 musculoskeletal shoulder examination using the PhyES. A kappa coefficient of $\kappa = 0.79$ was
271 reported.

272 Swift et al. (2013) published an ICC of 0.77 for the interrater reliability of the mO-S3 based
273 on the clinical competency of doctoral physical therapy students halfway through their
274 education in musculoskeletal physiotherapy.

275 An interrater reliability of ICC = 0.95 for students scored on all 138 items on the head to toe
276 examination (HTTPE) was reported by Yudkowsky et al. (2004). Lastly, Boulet et al. (2004)
277 reported a correlation coefficient of $r = 0.83$ (range $r = 0.06 - r = 0.93$) for the interrater
278 reliability of the OMT. The authors reported that the average difference between two raters
279 was 2.4 points on a 0 to 30 points scale. All interrater reliability estimates are presented in
280 Figure 3.

281

282 **Figure 3. Interrater reliability estimates.**

283

284 Intrarater reliability was available for only one assessment. Ladyshevsky et al. (2000)
285 published an intrarater reliability of $\kappa = 0.63$ for the PhyES.

286 None of the studies included in this review evaluated the measurement error of their
287 included assessments.

288 **4.2.2 Validity**

289 **Validity of the included assessments was evaluated with regard to their content validity,**
290 **criterion validity and construct validity.**

291 **Evidence for content validity was found for four assessments AMPE, PhyES, GPSE and mO-S3**
292 **(Beran et al., 2012; Ladyshevsky et al., 2000; Nothnagle et al., 2010; Swift et al., 2013).** For
293 each assessment, the authors provided information about how their assessments were
294 designed. All four studies used expert panels to judge comprehensiveness and relevance of
295 the assessment items. The size of the expert panels ranged between an unspecified numbers
296 of **panel members for the AMPE and mO-S3 (Beran et al., 2012; Swift et al., 2013) to 17**
297 **participants for the GPSE (Nothnagle et al., 2010).** Additionally, **two** studies involved learners
298 in the process of designing the assessment PhyES and GPSE (Ladyshevsky et al., 2000;
299 Nothnagle et al., 2010) with Nothnagle et al. (2010) generating content for the GPSE through
300 three focus groups. **None of the studies within this review reported the criterion validity of**

301 their assessments. Therefore, the utility of using the assessments to predict future
302 performance or compare to another measure is not known.

303 Data regarding the construct validity was available for five assessments AMPE, OMT, PES-C,
304 PES-R, PhyES (Beran et al., 2012; Boulet et al., 2004; Herbers et al., 2003; Ladyshevsky et al.,
305 2000). Three studies tested hypotheses whether their assessments could discriminate
306 performance between individuals with more experience or less experience. Beran et al.
307 (2012) reported that years of training had no significant influence on the total score of the
308 AMPE. Ladyshevsky and colleagues found in their study that licenced physiotherapists
309 performed significantly better on the PhyES than fourth year undergraduate students. Lastly,
310 Herbers et al. (2003) presented evidence, that learners in a training group scored
311 significantly higher than learners without a specific training ($p < 0.001$) on the PES-C. Two
312 studies reported correlations between the included assessments and other, established
313 assessments as evidence for construct validity. Herbers et al. (2003) reported an agreement
314 of $K = 0.66$ between their checklist for a pelvic examination (PES-C) and a global rating scale
315 for this procedure (PES-R). Boulet et al. (2004) reported that the OMT instrument correlated
316 with biomedical knowledge indicators ($r = 0.47$) and global patient assessment ($r = 0.46$).

317 4.3 Methodological quality of assessments

318 Methodological quality of the included assessments was low to moderate. Methodological
319 quality was appraised with 20 standards of the SEQAM. The assessment that was appraised
320 as fulfilling the most standards was the AMPE. Ten of the 20 standards were appraised as
321 fulfilled. The mO-S3 was evaluated as fulfilling the least standards (7 standards were
322 classified as satisfied). All standards are presented in Table 4.

323 Table 4. Methodological quality of included assessments

324 5 Discussion

325 The discussion is divided into the following sections: 1) summary of main results 2)
326 methodological quality of the assessments 3) potential biases in the review process and 4)
327 agreements and disagreements with other studies.

328 5.1 Summary of main results

329 This systematic review synthesised relevant literature relating to the current knowledge of
330 assessments for procedural skills in physiotherapy education. Following a systematic search,
331 eight assessments for procedural skills were identified that can be used in physiotherapy
332 education. Six of the assessments were designed for a specific procedure and were validated
333 for diagnostic or evaluative procedures. Two assessments (GPSE and OMT) were considered
334 useful for the evaluation of more than one procedure and can be used to evaluate
335 procedural competence of therapeutic interventions.

336 The GPSE was classified as representing the highest level of Miller's framework of
337 assessments (Miller, 1990) and can be used as workplace based assessment, which is the
338 "Does" level in Miller's pyramid. All the remaining assessments were classified as
339 representing the "Shows how" level, because they were all based in a simulated
340 environment and no direct evidence was available to evaluate whether the behaviour of the
341 learners actually changed.

342 In terms of internal consistency, the best performing assessment, (OMT), had a value above
343 0.70, while the other assessment reporting internal consistency (the mO-S3) had lower
344 estimates. These lower values of the mO-S3 might be explained by the method to calculate
345 internal consistency which was used by Swift et al. (2013). They calculated internal
346 consistency with regard to a 6 station OSCE, with stations designed to measure competence
347 in musculoskeletal physiotherapy. However, the content of the stations varied to some
348 extent. This conflicts with the stance of Cortina (1993) who stated that when internal
349 consistency is measured, the set of test items should form a reflective model, i.e. that "all
350 items are a manifestation of the underlying construct" (Mokkink et al., 2009, p. 24). It could
351 be argued that the stations and test items of the OSCE devised by Swift et al. (2013) did not
352 measure the same construct (e.g. diagnostic, interventional or communication competence)

353 or that they measured different aspects of one construct. This could explain the lower
354 internal consistency estimates of the mO-S3.

355 Six of the included assessments reported interrater reliability. The highest estimate was
356 reported for the HTTPE (ICC: 0.95). The AMPE and the PES-C were evaluated as having
357 moderate to low interrater reliability because estimates were below 0.70. There are a
358 number of methodological issues that may have affected the reliability. For the PES-C,
359 Herbers et al. (2003) calculated their reliability scores based on a subset of their items (i.e.
360 only data of 7 of the 29 items of the PES-C were used). Additionally, the study used
361 audiotapes to calculate the reliability between two raters. With regard to a checklist that
362 aims to evaluate procedural skills important issues may have been missed, which can only be
363 detected visually. Therefore, only such items as: "Asks if patient wants mirror to watch
364 examination" were evaluated with regard to their reliability. In relation to the AMPE, three
365 out of the total of four different assessments scored around or above the 0.7-margin. Only
366 the AMPE assessment of a physical examination of trauma patients scored considerably
367 lower (ICC = 0.27). Beran et al. (2012) reported that considerable disagreement was present
368 between raters. Especially, one rater scored consistently higher than the two other raters. In
369 an attempt to improve the reliability, the scores of three raters were averaged and
370 compared with an external rating. This method resulted in increased interrater reliability
371 scores (ICC = 0.51).

372 Only the PhyES evaluated intrarater reliability, reporting a moderate agreement ($\kappa = 0.63$).
373 These findings should be interpreted with caution due to the very small sample (six
374 encounters over two occasions during a two-week period).

375 When a new assessment is developed, users require reassurance that the instrument is
376 comprehensive and relevant. This might be assured by using experts to comment on or
377 generate the content of the assessment (Mokkink et al., 2009). Furthermore, the proposed
378 assessment should match the target population with regard to focus and detail and one way
379 of assuring this is to recruit potential participants and discuss the assessment with them.
380 However, only the PhyES (Ladyshevsky et al., 2000) and the GPSE (Nothnagle et al., 2010)
381 included students into the design of the assessments. Nothnagle et al. (2010) also used a
382 more robust development process, including focus groups, to construct their assessment

383 (GPSE), which may make it more likely that this assessment is comprehensive and consists of
384 relevant items.

385 Evidence of construct validity was found for four assessments (PES-C, PES-R, PhyES and
386 OMT). It has been established that learners should improve execution of a procedure in
387 response to level of experience and increased amounts of practice (Brydges, Carnahan,
388 Backstein, & Dubrowski, 2007). Specifically, the PES-C and the PhyES were able to
389 differentiate between learners with different levels of experience, however this was not
390 established for the AMPE.

391

392 5.2 Methodological quality of assessments

393 Methodological quality of assessments was evaluated with the SEQAM, which is based on
394 the utility index of Van Der Vleuten (1996). The author argued that appraisal of assessment
395 methods in health professions education should consider more than traditional
396 measurement properties (i.e. reliability and validity). Within his utility index he also placed
397 weight on the acceptability, the educational impact and the cost effectiveness of an
398 assessment. Because educators should take this information into account when context
399 specific decisions about assessments are made (Van Der Vleuten & Schuwirth, 2005). Similar
400 the SEQAM critically evaluates six dimensions: reliability, validity, ease of use, resources
401 required, ease of interpretation and educational impact.

402 Overall the methodological quality of the included assessments was low to moderate
403 (fulfilling between 6 and 10 standards). No assessment was appraised as having no risk of
404 bias. No study fulfilled all educational standards of the SEQAM. The assessment that was
405 appraised as fulfilling the most standards was the AMPE with 10 of the 20 standards fulfilled.
406 The mO-S3 was evaluated as fulfilling the least standards (6/20). The remaining assessments
407 ranged in between seven – nine standards fulfilled. One reason for this moderate quality of
408 evidence was that it was derived from only a single study for each assessment. Therefore, it
409 was not possible to complete some standards (e.g. the item “positively affects programme
410 curriculum” can only be awarded if at least two studies present evidence).

411 A discrepancy existed between the assessment and the standard “training requirements”.

412 The standard sets the benchmark for the training time to one hour in order to reduce the

413 required resources. In contrast, most of the researchers spent considerable more time in the
414 training of faculty members and standardised patients, with Ladyshevsky et al. (2000)
415 spending up to 30 hours in the training of their assessors. This is not viable in an educational
416 programme and therefore finding a reasonable balance between those extremes will be a
417 challenge for further work.

418 Within the “non-traditional” categories of measurement properties, (e.g. non- psychometric
419 properties) it was noted that five assessments were classified as “relatively easy to use”
420 because they required little specialist set up and time to evaluate (Beran et al., 2012; Boulet
421 et al., 2004; Nothnagle et al., 2010; Swift et al., 2013; Yudkowsky et al., 2004). However, only
422 the GPSE was as appraised as also requiring few resources (Nothnagle et al., 2010). This
423 could be important for educators when they need assessments in their daily practice which
424 are easy to set up and use.

425 5.3 Potential biases in the review process

426 Only one study for each assessment was identified, therefore limiting generalisability and
427 rendering it impossible to perform a meta-analysis. Findings have therefore been presented
428 narratively. Furthermore, sample size may affect findings, only three studies evaluated their
429 assessments with considerable sample sizes. Boulet et al. (2004), Herbers et al. (2003), and
430 Yudkowsky et al. (2004) used at least 70 participants in their studies. The remaining studies
431 recruited considerably fewer (< 25) participants which again may limit generalisability and
432 may have caused imprecision of the effect estimates.

433 A cut off value of 0.7 was used for the measurement properties of internal consistency and
434 interrater reliability and intrarater reliability (Terwee et al., 2007). While other authors use
435 different cut off values – e.g. 0.85 cut off (Weiner and Stewart (1998), the more moderate
436 interpretation was selected as 0.85 may be too high to be useful in practical settings
437 (Streiner, Norman, & Cairney, 2014). An acceptable reliability standard should be chosen
438 with regard to a specific situation. In high stakes examinations (i.e. tests with serious
439 consequences for the tester in situations such as education or certification (Sackett, Schmitt,
440 Ellingson, & Kabin, 2001)) higher reliability is required compared to a low stakes
441 examinations (i.e. tests without serious consequences for the learner).

442 A further potential bias in this review is that the SEQAM grading of the methodological
443 quality of assessment was modified. Swing et al. (2009) originally suggested an overall

444 recommendation (i.e. class of evidence) based on the evidence levels provided for each
445 dimension. We decided against the use of an overall score because firstly, in our view, scores
446 should only be combined when they are unidimensional (i.e. the same attribute of the object
447 “methodological quality” should be measured with the different sub-categories) and
448 evidence for unidimensionality was not available for SEQAM; secondly, the use of summary
449 scores might lead to biased estimates in systematic reviews and meta-analysis (da Costa,
450 Hilfiker, & Egger, 2013; Juni, Altman, & Egger, 2001). Therefore, we decided to omit the
451 overall recommendations and present relevant methodological aspects individually.

452 **5.4 Agreements and disagreements with other studies or reviews**

453 Four recent systematic reviews were identified that reported assessment of procedural skills
454 in health professions education (Bould, Crabtree, & Naik, 2009; Jelovsek et al., 2013;
455 McKinley et al., 2008; Morris et al., 2012).

456 In general, these reviews focussed on medical education and few assessments relevant for
457 use by allied health professions were identified. For example, of the assessments evaluated
458 in this review, only the OMT scale was identified by McKinley and colleagues. The remaining
459 assessments were not discussed in other reviews. Existing reviews do however agree that
460 there is a lack of assessments for procedural skills in allied health profession. In contrast, a
461 considerably greater number of assessments are available for use in medical education:
462 McKinley et al. (2008) included 85 different scales in their review of assessments used in
463 medical education. Our findings were similar to those of Jelovsek et al. (2013) who found
464 that there was limited reporting of measurement properties. Bould et al. (2009) suggested
465 that procedure unspecific assessments tended to miss errors in safety issues. We were not
466 able to comment as only two procedure unspecific assessments were included in this review
467 and this is therefore an area where uncertainty remains and further work is required.

468 **6 Conclusion and implications**

469 Following this systematic review, it was not possible to recommend a single assessment of
470 procedural skills in physiotherapy education, all the assessments we identified have
471 elements of strength and weakness. Therefore, evaluators should use existing tools carefully
472 when evaluating procedural performance of physiotherapy students. Most assessments we
473 identified were developed for use within the speciality of musculoskeletal physiotherapy and

474 these could be integrated into educational practice. There is however, a need to develop
475 new assessments to allow valid and reliable assessments of the broader spectrum of
476 physiotherapeutic practice in other specialities (e.g neurological practice and respiratory
477 practice). When assessments are selected or developed, faculty members should carefully
478 consider issues such as the usefulness and possible interpretation of the findings as well as
479 the more well established focus on measurement properties such as validity and reliability.
480 This may help prevent neglect of issues of importance to relevant stakeholders. Future
481 studies aiming to design new assessments should involve all stakeholders in the design of
482 the content, use and scoring of the assessment. Furthermore, the construct(s) to be
483 measured should be clearly defined.

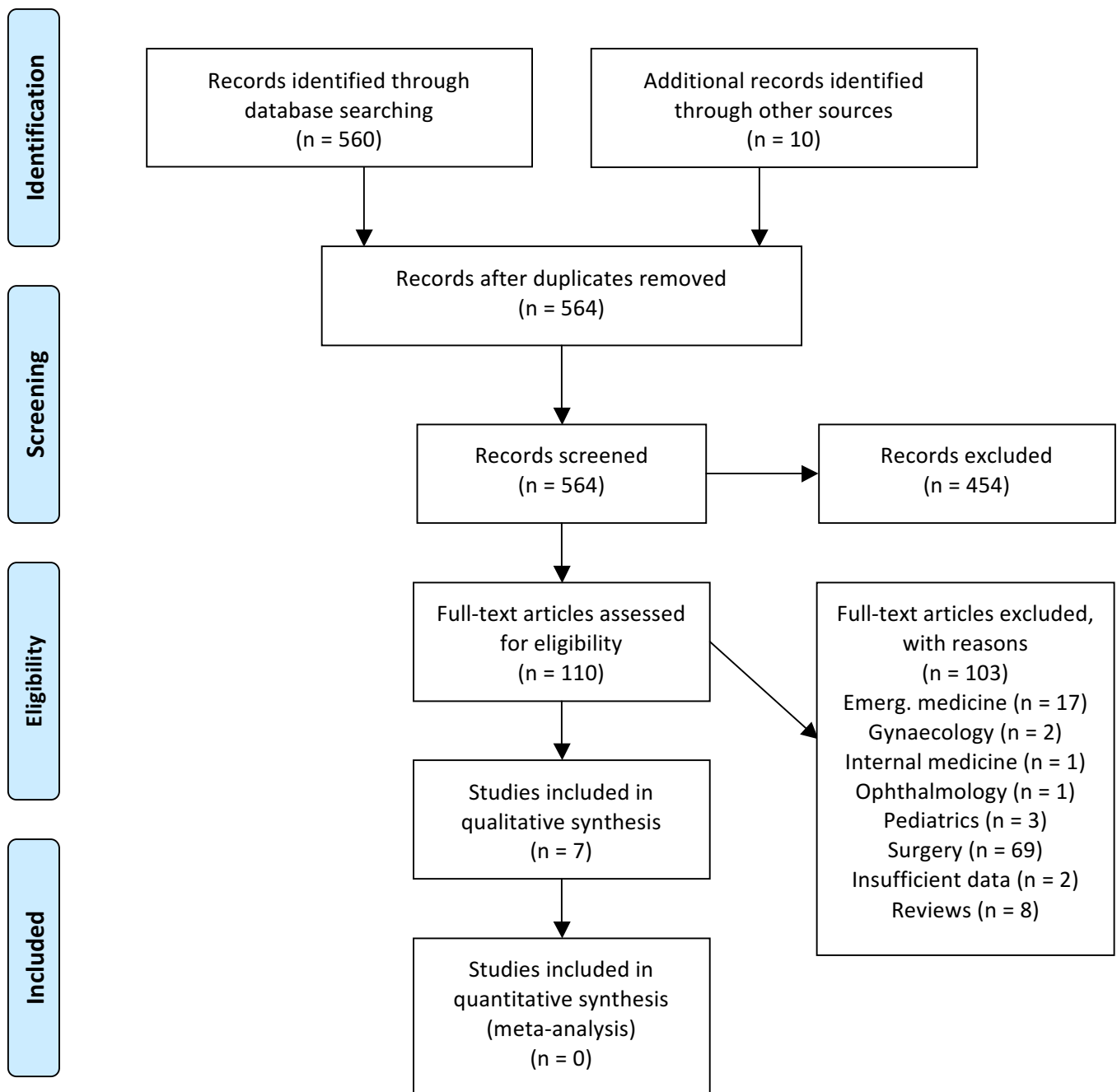
484

7 References

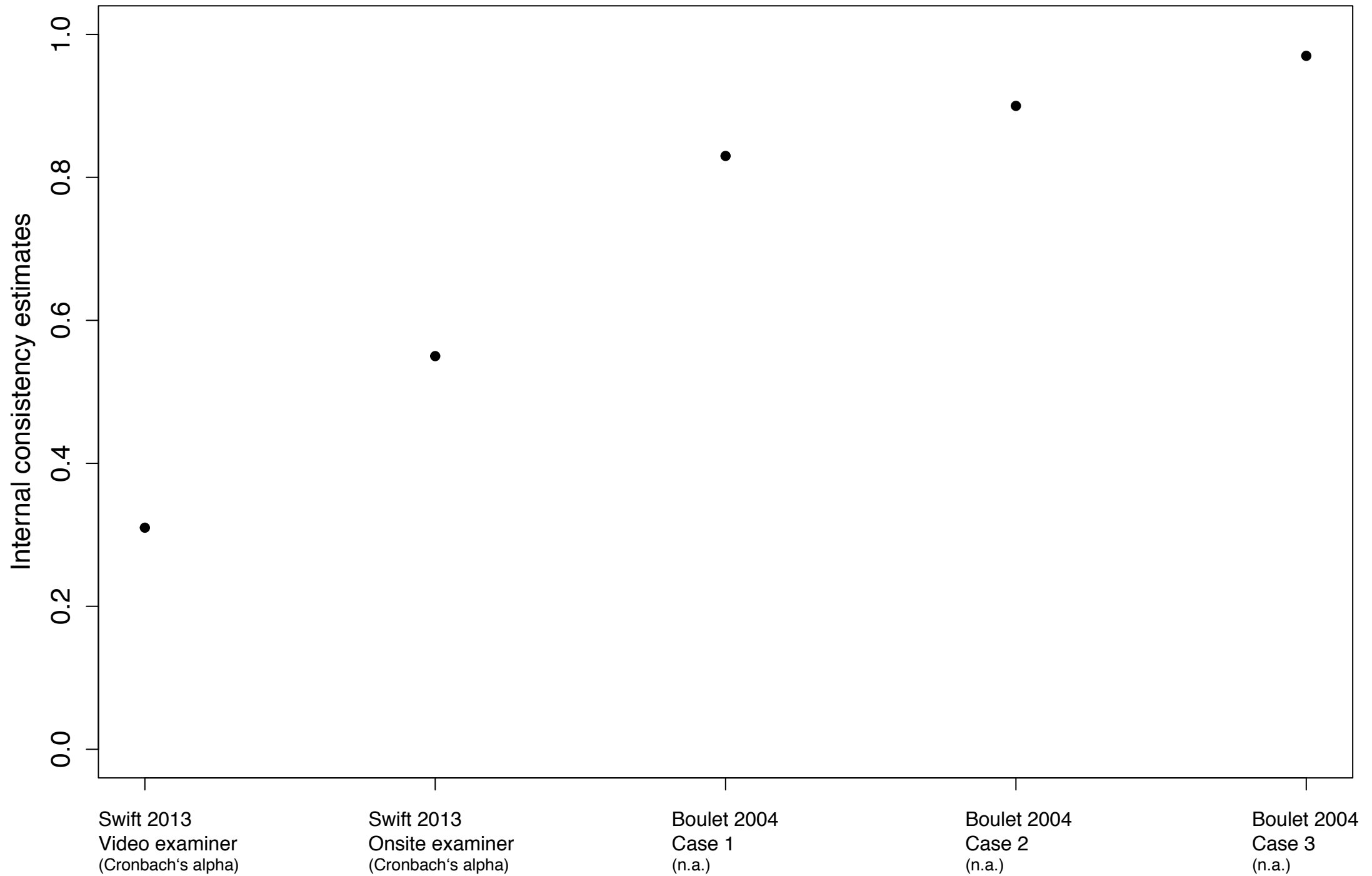
- Beran, M. C., Awan, H., Rowley, D., Samora, J. B., Griesser, M. J., & Bishop, J. Y. (2012). Assessment of musculoskeletal physical examination skills and attitudes of orthopaedic residents. *The Journal of Bone & Joint Surgery*, *94*(6), e36 31-38.
- Bould, M. D., Crabtree, N. A., & Naik, V. N. (2009). Assessment of procedural skills in anaesthesia. *Br J Anaesth*, *103*(4), 472-483.
- Boulet, J. R., Gimpel, J. R., Dowling, D. J., & Finley, M. (2004). Assessing the ability of medical students to perform osteopathic manipulative treatment techniques. *JAOA: Journal of the American Osteopathic Association*, *104*(5), 203-211.
- Brydges, R., Carnahan, H., Backstein, D., & Dubrowski, A. (2007). Application of motor learning principles to complex surgical tasks: searching for the optimal practice schedule. *J Mot Behav*, *39*(1), 40-48. doi:10.3200/jmbr.39.1.40-48
- Chandratilake, M., Davis, M., & Ponnampereuma, G. (2010). Evaluating and designing assessments for medical education: the utility formula. *Int J Med Educ*, *1*(1), 1-17.
- Cortina, J. M. (1993). What is coefficient alpha? An examination of theory and applications. *Journal of applied psychology*, *78*(1), 98.
- da Costa, B. R., Hilfiker, R., & Egger, M. (2013). PEDro's bias: summary quality scores should not be used in meta-analysis. *Journal of clinical epidemiology*, *66*(1), 75.
- General Medical Council. (2004). *The New Doctor: Guidance on PRHO training*. London: GMC.
- Glista, J., Pop, T., Weres, A., Czenczek-Lewandowska, E., Podgórska-Bednarz, J., Rykała, J., . . . Rusek, W. (2014). Change in anthropometric parameters of the posture of students of physiotherapy after three years of professional training. *BioMed research international*, 2014.
- Gorrell, L. M., Engel, R. M., Brown, B., & Lystad, R. P. (2016). The reporting of adverse events following spinal manipulation in randomized clinical trials—a systematic review. *The Spine Journal*.
- Herbers, J. E., Jr., Wessel, L., El-Bayoumi, J., Hassan, S. N., & St Onge, J. E. (2003). Pelvic examination training for interns: a randomized controlled trial. *Acad Med*, *78*(11), 1164-1169.
- Higgs, J., Hunt, A., Higgs, C., & Neubauer, D. (1999). Physiotherapy education in the changing international healthcare and educational contexts. *Advances in Physiotherapy*, *1*(1), 17-26.
- Jackson, J., & Liles, C. (1994). Working postures and physiotherapy students. *Physiotherapy*, *80*(7), 432-436.
- Jelovsek, J. E., Kow, N., & Diwadkar, G. B. (2013). Tools for the direct observation and assessment of psychomotor skills in medical trainees: a systematic review. *Med Educ*, *47*(7), 650-673.
- Juni, P., Altman, D. G., & Egger, M. (2001). Assessing the quality of controlled clinical trials. *British Medical Journal*, *323*(7303), 42.
- Kent, M. (2007). *The Oxford Dictionary of Sports Science & Medicine* (3 ed.). Oxford: Oxford University Press.

- Ladyshevsky, R., Baker, R., Jones, M., & Nelson, L. (2000). Evaluating clinical performance in physical therapy with simulated patients. *Journal of Physical Therapy Education*, 14(1), 31.
- Liberati, A., Altman, D. G., Tetzlaff, J., Mulrow, C., Gotzsche, P. C., Ioannidis, J. P., . . . Moher, D. (2009). The PRISMA statement for reporting systematic reviews and meta-analyses of studies that evaluate health care interventions: explanation and elaboration. *PLoS Medicine*, 6(7), e1000100. doi:10.1371/journal.pmed.1000100
- 10.1371/journal.pmed.1000100
- McKinley, R. K., Strand, J., Ward, L., Gray, T., Alun - Jones, T., & Miller, H. (2008). Checklists for assessment and certification of clinical procedural skills omit essential competencies: a systematic review. *Med Educ*, 42(4), 338-349.
- Michels, M. E., Evans, D. E., & Blok, G. A. (2012). What is a clinical skill? Searching for order in chaos through a modified Delphi process. *Med Teach*, 34(8), e573-e581.
- Miller, G. E. (1990). The assessment of clinical skills/competence/performance. *Academic medicine*, 65(9), S63-67.
- Mokkink, L. B., Terwee, C. B., Patrick, D. L., Alonso, J., Stratford, P. W., Knol, D. L., . . . De Vet, H. C. (2009). The COSMIN checklist manual. *Amsterdam: VU University Medical Centre*.
- Mokkink, L. B., Terwee, C. B., Patrick, D. L., Alonso, J., Stratford, P. W., Knol, D. L., . . . de Vet, H. C. (2010). The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. *Journal of clinical epidemiology*, 63(7), 737-745.
- Morris, M. C., Gallagher, T. K., & Ridgway, P. F. (2012). Tools used to assess medical students competence in procedural skills at the end of a primary medical degree: a systematic review. *Med Educ Online*, 17.
- Nothnagle, M., Reis, S., Goldman, R., & Diemers, A. (2010). Development of the GPSE: a tool to improve feedback on procedural skills in residency. *Fam Med*, 42(7), 507-513.
- Nyland, L. J., & Grimmer, K. A. (2003). Is undergraduate physiotherapy study a risk factor for low back pain? A prevalence study of LBP in physiotherapy students. *BMC musculoskeletal disorders*, 4(1), 22.
- Sackett, P. R., Schmitt, N., Ellingson, J. E., & Kabin, M. B. (2001). High-stakes testing in employment, credentialing, and higher education: Prospects in a post-affirmative-action world. *American Psychologist*, 56(4), 302.
- Simpson, J., Furnace, J., Crosby, J., Cumming, A., Evans, P., David, M. F. B., . . . McLachlan, J. (2002). The Scottish doctor--learning outcomes for the medical undergraduate in Scotland: a foundation for competent and reflective practitioners. *Med Teach*, 24(2), 136-143.
- Streiner, D. L., Norman, G. R., & Cairney, J. (2014). *Health measurement scales: a practical guide to their development and use*: Oxford university press.
- Swift, M., Spake, E., & Gajewski, B. J. (2013). The Reliability of a Musculoskeletal Objective Structured Clinical Examination in a Professional Physical Therapist Program. *Journal of Physical Therapy Education*, 27(2), 41.
- Swing, S. R., Clyman, S. G., Holmboe, E. S., & Williams, R. G. (2009). Advancing resident assessment in graduate medical education. *Journal of graduate medical education*, 1(2), 278-286.
- Terwee, C. B., Bot, S. D., de Boer, M. R., van der Windt, D. A., Knol, D. L., Dekker, J., . . . de Vet, H. C. (2007). Quality criteria were proposed for measurement properties of health status questionnaires. *Journal of clinical epidemiology*, 60(1), 34-42.

- The Royal Australian College of General Practitioners. (2011). *Procedural Skills*.
- Van Der Vleuten, C. P. (1996). The assessment of professional competence: developments, research and practical implications. *Advances in Health Sciences Education, 1*(1), 41-67.
- Van Der Vleuten, C. P., & Schuwirth, L. W. (2005). Assessing professional competence: from methods to programmes. *Med Educ, 39*(3), 309-317.
- Weiner, E. A., & Stewart, B. J. (1998). *Assessing individuals*: Brooks/Cole Publishing Company.
- Yudkowsky, R., Downing, S., Klamen, D., Valaski, M., Eulenberg, B., & Popa, M. (2004). Assessing the head-to-toe physical examination skills of medical students. *Med Teach, 26*(5), 415-419. doi:10.1080/01421590410001696452



Internal consistency



Interrater reliability

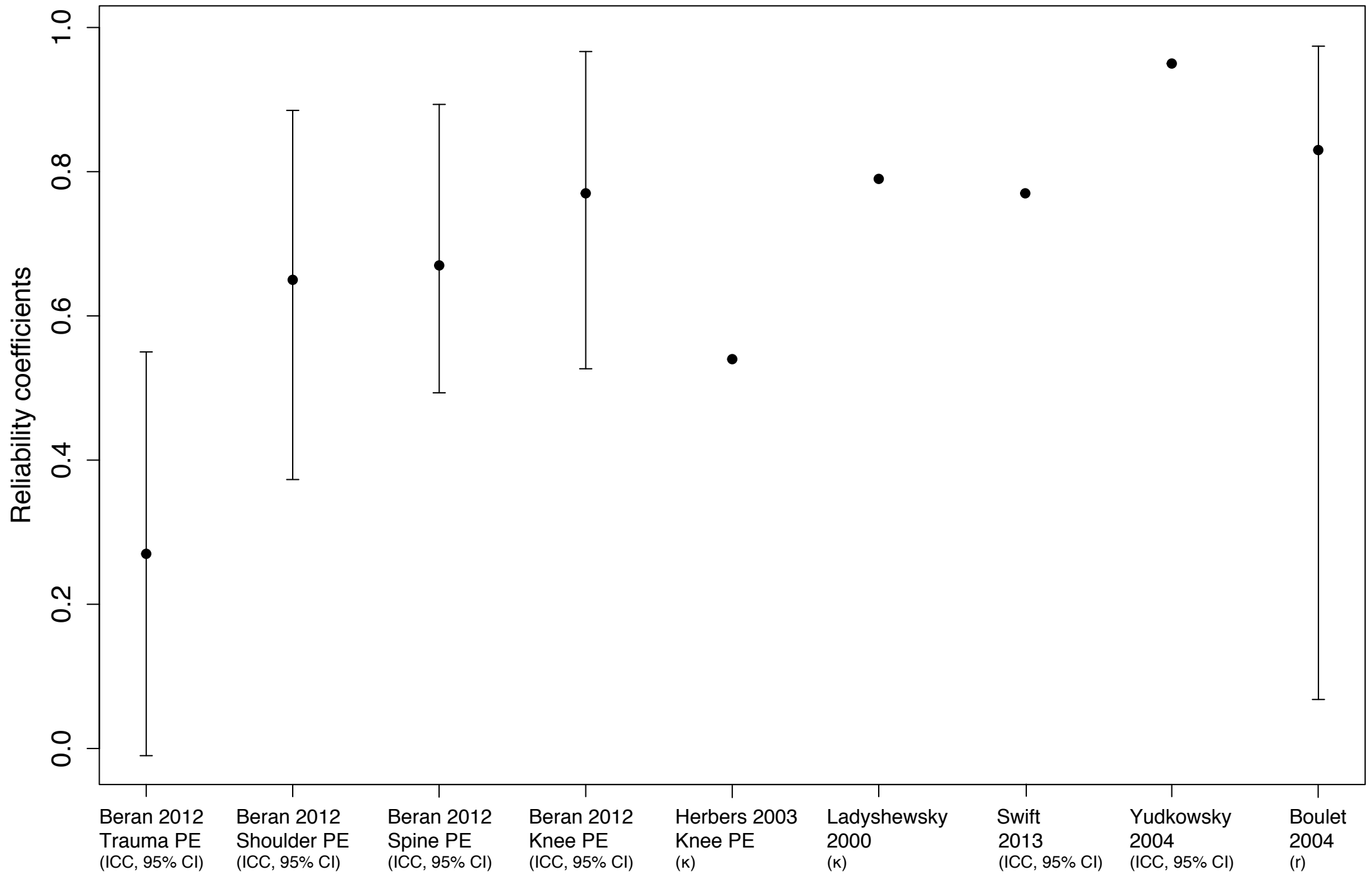


Table 1. In-and exclusion criteria

Category	Criteria
Population	Studies with physiotherapists or physiotherapy students were included.
	Studies with health professionals or health professions students were included when they practiced procedures that can be used in physiotherapy (i.e. when medical students were evaluated on their ability to perform a musculoskeletal examination)
	Studies with health professionals or health professions students were excluded when they practiced procedures that cannot be practiced by physiotherapists (such as surgery)
Educational assessments	The assessment could be either a procedure specific measurement instrument (i.e. the assessment is designed exclusively for one procedure) or a procedure unspecific measurement (i.e. the assessment is designed to measure procedures in physiotherapy education but can be used for more than one procedure)
	The assessment should measure procedures in reality. Assessments based on virtual reality were excluded.
	The assessment should be feasible in various settings. Therefore, assessments that require expensive equipment were excluded.
	Data must be available for a specific assessment. Studies with summary data of several assessments were excluded (e.g. summary scores of a complete OSCE).
Outcome	The aim of the assessment should be to measure procedural skills. Assessments of similar constructs such as clinical skills or psychomotor skills (defined as "... motor skill, some manipulation of material, or some act which requires a neuromuscular action" Simpson (1966, p. 17)) were included.
	Assessments that aimed to exclusively evaluate other outcomes such as communication skills or professionalism were excluded.
	When assessments were designed to measure multiple outcomes it was evaluated whether the main focus was based on procedural skills (e.g. more than 50% of the items concentrate on procedural skills). Assessments with the main focus on procedural skills were included.
Measurement properties	Studies had to report the measurement properties of an educational assessment (e.g. reliability or validity)

1. Reference

Simpson, E. J. (1966). The Classification of Educational Objectives, Psychomotor Domain.

Table 2. Search strategy

Population	Assessment	Outcome
"medical education" OR education, medical[Mesh] OR "physiotherapy education" OR "physical therapy education" OR "health professions education" OR "healthcare education" OR "allied health care education"	scale OR global rating scale OR GRS OR checklist	practical skill* OR psychomotor skill* OR procedural skill* OR clinical skill*

Table 3. Characteristics of included studies and assessments

Study	Country	Setting	Sample	Assessed procedure	Scale and items	Duration	Patients	Assessors	Purpose
Beran et al. (2012) AMPE	USA	Orthopaedic department	24 orthopaedic residents	PSMI: Musculoskeletal physical examination; Inspection, palpation, joint range of motion, strength testing and any special tests pertinent to the clinical scenario	Four 12-15 items checklists for clinical scenarios (upper extremity, lower extremity, trauma and spine) on dichotomous scales (yes or no).	10 Minutes	Standardised patients are required (120 minutes training)	Pool of experienced raters	High stakes purpose
Boulet et al. (2004) OMT	USA	Osteopathic college	121 osteopathic students (4 th year)	PUMI: Osteopathic manipulative treatment of three clinical cases (low back pain, frozen shoulder and asthmatic with cough)	OMT (Osteopathic Manipulative Treatment) assessment tool with 15 items; Every item is scored on a 0 to 2 scale (0 = done incorrectly or not done, 1 = not performed optimally and 2 = done proficiently)	13 minutes	Standardised patients with 8 hours of formal training	16 osteopathic physicians (5 hours of formal training)	High stakes examination (OSCE)
Herbers et al. (2003) PES-C	USA	University Medical Centre	72 internal medicine residents	PSMI: Pelvic examination	29 item dichotomous checklist (yes = when the behaviour was observed; no = when the behaviour was not observed); Includes some items about communication skills	Not specified	Gynaecologic teaching trainer required; 1 trainer was being examined and the second trainer rated the student's skills.	Gynaecologic teaching trainer required	Not specified

Herbers et al. (2003) PES-R	USA	University Medical Centre	72 internal medicine residents	PUMI: Pelvic examination	Global rating scale evaluating the overall performance of the pelvic examination (five-point ordinal scale between 1 = inadequate and 5 = excellent)	Not specified	Gynaecologic teaching trainer required; 1 trainer was being examined and the second trainer rated the student's skills.	Gynaecologic teaching trainer required	Not specified
Ladyshevsky et al. (2000) PhyES	Australia	Physiotherapy department	12 undergraduate physiotherapy students 4 physiotherapists (at least 2 years of experience)	PSMI: Musculoskeletal physical examination of a patient with a rotator cuff problem	Physical examination checklist (3 point scale: 0 = not done, 1 = done poorly or incompletely and 2 = done well), number of items not available	Mean 30 minutes (range: 20 - 46 minutes)	Standardised patients are required	Assessors with 30 hours of training	High stakes examination (OSCE)
Nothnagle et al. (2010) GPSE	USA	Family medicine department	5 faculty members and 5 students (semi structured interviews); Focus groups: 7 experienced family medicine educators, 5	PUMI: Eligible for all procedures in family medicine	Global Procedural Skills Evaluation Form, 4 point scale, amount of assistance is documented ranging from significant guidance is provided to performed independently; communication skills etc. are included; Student's self assessment is included;	Not specified	Not required	Not specified	Low stakes examination (formative feedback)

			residents and 5 faculty members		Difficulty of the procedures is rated as well				
Swift et al. (2013)* mO-S3	USA	Physiotherapy department	12 undergraduate 1 st year physiotherapy students	PSMI: Examination skills in musculoskeletal physiotherapy (shoulder tests)	Checklist for a musculoskeletal OSCE station; 6 item checklist (5 dichotomous items and 1 ordinal item)	6 minutes	Simulation patients with 2 hours of supervised training and 1 week of independent training	Clinical instructors (2 - 20 years of experience)	Low stakes examination (mid-term)
Yudkowsky et al. (2004) HTTPE	USA	University Medical Centre	369 medical students (2 nd year)	PSMI: Head to toe physical examination	138 item checklist; three-point scale (0 = incorrect, 1 = correct after prompt, 2 = correct without prompting); Test duration: 2h;; high stakes summative assessment or low stakes formative assessment	2 hours (45 minutes unprompted exam, remaining 1:45 hours are used for scoring, feedback, and teaching)	Trained patient instructors with 25 hours of training	Trained patient instructors with 25 hours of training	High stakes summative assessment and low stakes formative assessment

AMPE: Assessment of Musculoskeletal Physical Examination Skills Checklist; GPSE: Global Procedural Skills Evaluation Form; HTTPE: Head to Toe Physical Examination; mO-S3: mOSCE-Station 3 checklist; PhyES: Physical Examination Skills Checklist; PES-C: Pelvic Examination Skills Checklist; PES-R: Pelvic Examination Skill Rating Scale; PSMI: Procedure Specific Measurement Instrument; PUMI: Procedure Unspecific Measurement Instrument

* Swift et al. (2013): It was only possible to use data from a small pilot study. The follow up study evaluated a 6 station OSCE. Single values for a specific scale were not available.

Table 4. Methodological quality of included assessments

Standards for evaluating the quality of assessment methods	Beran 2012 (AMPE)	Boulet 2004 (OMT)	Herbers 2003 (PES-C&R)	Ladyshewsky 2000 (PhyES)	Nothnagle 2010 (GPSE)	Swift 2013 (mO-S3)	Yudkowsky 2004 (HTTPE)
Reliability							
1. Reliability indicators	☺	☺	☹	☺	☹	☺	☺
2. Inter- and Intrarater reliability	☹	☹	☹	☺	☹	☹	☹
3. High-stakes decisions	☹	☹	☹	☹	☹	-	☹
Level of evidence (A, B, C or NR)	C	C	NR	B	NR	C	C
Validity							
1. Interpretation of results	☺	☺	☺	☺	☺	☹	☺
2. Selection of content	☺	☹	☺	☺	☺	☺	☺
3. Unintended consequences	☺	☺	☺	☺	☺	☺	☺
4. Agreement between a single expert and consensus ratings	☺	☹	☹	☹	☹	☹	☹
5. Subjective judgment	☹	☹	☹	☹	☹	☹	☹
Level of evidence (A, B, C or NR)	B	NR	NR	C	NR	NR	NR
Ease of use							
1. Daily practice	☹	☹	☺	☹	☺	☹	☺
2. Special set up	☺	☺	☺	☺	☺	☺	☺
3. Duration	☺	☺	☺	☹	☺	☺	☹
Level of evidence (A, B, C or NR)	B	B	C	C	B	B	B
Resources required							
1. Additional resources	☺	☺	☺	☺	☺	☺	☺
2. Training requirements	☹	☹	☹	☹	☺	☹	☹
3. Additional persons	☹	☹	☹	☹	☺	☹	☹
Level of evidence (A, B, C or NR)	C	C	C	C	B	C	C
Ease of interpretation							
1. Interpretation of individual scores	☺	☺	☺	☺	☺	☹	☺
2. Normative data	☺	☹	☺	☹	☹	☹	☹
3. Individual to group performance.	☹	☹	☹	☹	☹	☹	☹
Level of evidence (A, B, C or NR)	B	C	B	NR	C	NR	C
Educational impact							
1. Positively affect individual learners	☹	☹	☺	☹	☹	☹	☺
2. Positively affect programme curriculum	☹	☹	☹	☹	☹	☹	☹
3. Provide useful results	☺	☺	☺	☺	☺	☺	☺
Level of evidence (A, B, C or NR)	NR	NR	C	NR	C	NR	B

A level of evidence A was assigned when all standards in one dimension were met. A level of B was assigned when one standard was not met. A level of C was appraised when two standards were not met and NR was assigned when more than two standards were not met. ☹: Standard not met; ☺: Standard met; ☺: Unclear; -: Standard not applicable