# 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 the Global Rating of Change Scales in Patients with Neck Disorders: A Systematic Review with Meta-Analysis and Meta-Regression

SCHOLARONE™
Manuscripts

1
2
3    **Psychometric Properties of the Global Rating of Change Scales in Patients with Neck**
4
5    **Disorders: A Systematic Review with Meta-Analysis and Meta-Regression**
6
7
8    Pavlos Bobos[1], Joy C MacDermid[2], Goris Nazari[3], Rochelle Furtado[4] and CATWAD co-authors[5]
9
10   4
11
12
13
14   [1]Pavlos Bobos PT, PhD(c), (corresponding author) Doctoral Candidate, Western's Bone and Joint
15   Institute, Department of Health and Rehabilitation Sciences, Western University, Elborn College,
16   1201 Western Road, N6G 1H1, London, Ontario, Dalla Lana School of Public Health, Institute of
17   Health Policy Management and Evaluation, Department of Clinical Epidemiology and Health Care
18   Research, University of Toronto, Canada, (pbobos@uwo.ca), tel: +1 519 661 2111 x88912

19   [2]Joy C MacDermid BScPT, PhD, Professor, Physical Therapy and Surgery, Western University,
20   London, ON and Co-director Clinical Research Lab, Hand and Upper Limb Centre, St. Joseph's
21   Health Centre, London, Ontario; Professor Rehabilitation Science McMaster University,
22   Hamilton, ON, Canada (jmacderm@uwo.ca)

23   [3]Goris Nazari PT, PhD(c) Doctoral Candidate, Western's Bone and Joint Institute, School of
24   Physical Therapy, Department of Health and Rehabilitation Sciences, Western University,
25   London, Ontario, Canada, (gnazari@uwo.ca)

26   [4]Rochelle Furtado MSc Western's Bone and Joint Institute, School of Physical Therapy,
27   Department of Health and Rehabilitation Sciences, Western University, London, Ontario, Canada,
28   (rfurtad5@uwo.ca)

29   [5]CATWAD: Michele Sterling m.sterling@uq.edu.au, Anne Söderlund anne.soderlund@mdh.se,
30   Michele Curatolo, curatolo@uw.edu, James M Elliott j-elliott@northwestern.edu, David Walton
31   dwalton5@uwo.ca, Helge Kasch helgkasc@rm.dk, Linda Carroll linda.carroll@ualberta.ca,
32   Hans Westergren Hans.Westergren@skane.se, Gwendolen Jull g.jull@uq.edu.au, Eva-Maj
33   Malmström eva-maj.malmstrom@med.lu.se, Luke B Connelly l.connelly@uq.edu.au, Joy C
34   MacDermid jmacderm@uwo.ca, Mandy Nielsen mandy.nielsen@griffith.edu.au, Pierre Côté
35   pierre.cote@uoit.ca, Tonny Elmose Andersen tandersen@health.sdu.dk, Trudy Rebbeck
36   trudy.rebbeck@sydney.edu.au, Annick Maujean a.maujean@uq.edu.au, Sarah Robins
37   s.robins1@uq.edu.au, Kenneth Chen k.chen8@uq.edu.au, Julia Treleaven j.treleaven@uq.edu.au

1

31  **ABSTRACT**

32  **Objective:** The purpose of this systematic review was to critically appraise and synthesize the

33  psychometric properties of Global Rating of Change (GROC) scales for assessment of patients

34  with neck pain.

35  **Design:** Systematic review

36  **Data sources:** A search was performed in 4 databases (MEDLINE, EMBASE, CINAHL,

37  SCOPUS) until February 2019.

38  **Data extraction and synthesis:** Eligible articles were appraised using Consensus-based Standards

39  for the selection of health Measurement Instruments (COSMIN) checklist and the Quality

40  Appraisal for Clinical Measurement Research Reports Evaluation Form.

41  **Results:** The search obtained 16 eligible studies and included in total 1533 patients with neck pain.

42  Test-retest reliability of Global Perceived Effect (GPE) was very high (Intra-class correlation

43  coefficient (ICC) = 0.80 to 0.92) for patients with whiplash. Pooled data of Pearson's r indicated

44  that GROC scores were moderately correlated with neck disability change scores (0.53, 95% CI:

45  0.47 to 0.59). Pooled data of Spearman's correlations indicated that GROC scores were moderately

46  correlated with neck disability change scores (0.56, 95% CI: 0.41 to 0.68).

47  **Conclusions:** This study found excellent quality evidence of very good to excellent test-retest

48  reliability of GPE for patients with Whiplash Associated Disorders. Evidence from very good-to-

49  excellent quality studies found that GROC scores are moderately correlated to an external criterion

50  patient-reported outcome (PROM) measure evaluated pre-post treatment in patients with neck

51  pain. No studies were found that addressed the optimal form of GROC scales for patients with

52  neck disorders or compared the GROC to other options for single-item global assessment.

53  **Prospero registration number:** CRD 42018117874

54

2

**Strengths and limitations of this study**

- We rated the quality of individual studies and the overall risk of bias using two standardized approaches

- Our focus on neck pain increased the specificity of results but are not necessarily applicable to other musculoskeletal conditions

- Conceptual concerns about global ratings of change being affected by recall bias are not adequately addressed by psychometric evidence

- No studies addressing the optimal form of global rating were found.

**Introduction**

Neck pain is the 4th leading cause of disability and approximately half of adult the population with neck pain will experience a clinically important episode once in their lifetime. [1–3] The annual prevalence of neck pain it is estimated between 15% and 50%, with females having a higher prevalence rate than males. [2,3] Neck pain has been associated with many other comorbidities such as headaches, dizziness, anxiety, depression, back pain and arthralgias.[3–6] Several different methods for classifying neck pain have been described, using indicators such as duration (acute, sub-acute or chronic), degree of interference (low, moderate, severe) or most likely structure at fault (e.g. neuropathy vs. mechanical). [7]

As part of a patient-centric approach to care, clinicians will commonly evaluate response to intervention by asking the patient directly whether they feel better, worse, or the same since the prior encounter. While direct questioning can provide a qualitative indicator of change in status, many best practice guidelines endorse use of some form of quantified patient-reported outcome (PRO) as an adjunct to oral self-report. PROs are available to quantify several different constructs in people with neck pain, including pain severity, disability and neck function. [8] Any PRO

3

80   intended to provide an estimate of change over time should be responsive to subtle shifts in the

81   patient's condition. To facilitate interpretation of change scores, a common property of many such

82   tools is the minimum clinically important difference (MCID), which is a change threshold that

83   corresponds to the minimum shift in scale values that most patients would indicate corresponds to

84   an important change in their overall condition. A well-recognized approach to establishing an

85   MCID for a PRO is to compare the magnitude of change against an anchor, most commonly a

86   Global Rating of Change (GROC) scale. These scales allow patients or study participants to

87   indicate whether their condition has gotten worse, better, or stayed the same and to quantify the

88   magnitude of that change. As they have been adopted as a sort of 'standard' against which change

89   in other tools is compared, the GROC can also be used on its own as an omnibus generic indicator

90   of change. [8]

91   Despite being accepted as a standard measure, there is considerable variation in how the

92   GROC has been constructed and implemented in research in neck pain. Some are 15 points, some

93   11 points, and others are 7 points. The common structure across these is the use of a middle '0'

94   score corresponding to 'no change', with negative values indicating magnitudes of worsening

95   while positive values indicate improvement.[9] Variations of the GROC (in name or structure)

96   include the "Global Perceived Effect", "Patient Global Impression of Change", "Transition

97   Ratings", and "Global Scale". [9]

98   A critical component of monitoring changes in health outcomes is having valid, reliable

99   and responsive tools with strong psychometric properties. While recent research [8] has examined

100   the psychometric properties of the most commonly reported PROs for neck disorders, to date there

101   has been no systematic review to summarize the measurement properties of GROC scales

102   themselves in patients with neck disorders. Therefore, this systematic review aims to critically

4

103    appraise and synthesize the psychometric properties of the GROC scales in patients with neck

104    disorders.

105

106    **METHODS**

107    *Patient and Public Involvement*

108    There was no patient or public involvement in the design or planning of this study.

109

110    *Study Design and Protocol Registration*

111    We conducted a systematic review to evaluate the psychometric properties of GROC scales in

112    patients with neck disorders. The protocol was registered in PROSPERO register database with

113    registration number: CRD 42018117874

114

115    *Eligibility Criteria*

116    We included studies in this systematic review if the following criteria were met [10–12]:

- Design: psychometric testing, randomized/ cohort studies

- Participants: > 50% of the study's patient population with neck conditions/disorders,

- Intervention/Comparison: studies that reported on the psychometric properties (reliability, validity, responsiveness) of GROC, Global Perceived Effect (GPE) and Patient Global Impression of Change (PGIC),

- Outcomes: GROC, GPE and PGIC.

123    Studies with no data on the GROC scales' psychometric properties, and conference

124    abstract/posters were excluded from this systematic review.

125

5

126 *Information Sources*

127 To identify studies on the psychometric properties (reliability, validity, responsiveness) of the

128 GROC, GPE and PGIC we searched the Medline, EMBASE, Scopus and CINAHL databases from

129 inception till February 2019, using a combination of keywords. Furthermore, we identified

130 additional studies by examining the reference list of each of the selected studies. The full list with

131 keyword strategy is presented in **APPENDIX 1**.

132

133 *Study Selection*

134 Two investigators (PB and GN) performed the systematic electronic searches independently in

135 each database. The same investigators then proceeded to identify and remove the duplicate studies.

136 In the next stage, we performed the independent screening of the titles and abstracts and any full-

137 text article marked as include or uncertain were obtained. In the final stage, the same two

138 independent authors performed the full text reviews independently to assess final article eligibility.

139 In case of disagreement, a third reviewer; the most experienced member (JM), facilitated a

140 consensus through discussion.

141

142 *Data Extraction*

143 The fourth author (RF) performed the data extractions. The extracted data were then cross-checked

144 by another author (PB). Data extraction included the author, year, study population/condition,

145 setting, sample size, age, properties evaluated, retest-interval, and the intervention protocol (if used

146 to assess responsiveness parameters). [13,14] For reliability estimates, Standard Error of

147 Measurement (SEM), Intra-class Correlation Coefficient (ICC), Minimal Detectable Change

148 (MDC) and 95% confidence intervals were extracted. [13,14] The ICC interpretation of ICC < 0.40

149 indicating poor, $0.40 \leq \text{ICC} < 0.75$ indicating fair-to-good and $\text{ICC} \geq 0.75$ indicating excellent

150 reliability were used as a common benchmark. For validity estimates, correlation coefficient

151 (Pearson's/Spearman) and the 95% confidence intervals were extracted. [13,14] Evan's guidelines

152 to interpret the strength of the correlation was used which included: 0.00–0.19 "very weak", 0.20–

153 0.39 "weak", 0.40–0.59 "moderate", 0.60–0.79 "strong", and 0.80–1.00 "very strong". [15] For

154 responsiveness estimates, the Effect Size (ES), Standardized Response Mean (SRM), Clinically

155 Important Difference (CID), and/or Minimal Clinically Important Difference (MCID) including

156 the method of MCID estimation − Anchor-/Distribution-based methods, and 95% confidence

157 intervals were extracted. [13,14] To assist clinical decision making, standard benchmark scores of

158 trivial (< 0.20), small (≥ 0.20 to < 0.50), moderate (≥ 0.50 to < 0.80) or large (≥ 0.80), as proposed

159 by Cohen, were used. [16] When insufficient data were presented, PB contacted the authors by

160 email and requested further data.

161

162 *Consensus-based Standards for the selection of health Measurement Instruments (COSMIN)*

163 Consensus-based Standards for the selection of health Measurement Instruments (COSMIN)

164 assesses the risk of bias for the psychometric properties reported on a property-by-property basis.

165 A score for the risk of bias in estimates of psychometric properties was assessed by two authors

166 (PB) and (RF) using the new (COSMIN) checklist.[17] If disagreement was present a third person

167 (JM) assist in resolving the discrepancy. Each study was scored on the 4-point scale as "very

168 good", "adequate", "doubtful" or "inadequate" for each of the checklist criteria for relevant

169 measurement properties (e.g. reliability, responsiveness, etc.). To determine the overall score for

170 each measurement property, the worst score counts method was used wherein the lowest score for

171 the checklist criteria of the relevant property was taken as the overall score. [18] We then assessed

7

172  the result of individual studies on a measurement property against the updated criteria for good

173  measurement properties. This involved the evaluation of results of included studies as either

174  sufficient (+), insufficient (–), or indeterminate (?). [17]

175

176  *Quality Appraisal for Clinical Measurement Research Reports Evaluation Form*

177  A summary score for the overall quality of individual studies was appraised independently by the

178  authors (PB) and (RF) using a structured clinical measurement specific appraisal tool. [13,14] In

179  case of disagreement a third person was consulted (JM) to resolve the conflict. The evaluation

180  criteria of this tool included twelve items: 1) Thorough literature review to define the research

181  question; 2) Specific inclusion/exclusion criteria; 3) Specific hypotheses; 4) Appropriate scope of

182  psychometric properties; 5) Sample size; 6) Follow-up; 7) The authors referenced specific

183  procedures for administration, scoring, and interpretation of procedures; 8) Measurement

184  techniques were standardized; 9) Data were presented for each hypothesis; 10) Appropriate

185  statistics-point estimates; 11) Appropriate statistical error estimates; and 12) Valid conclusions

186  and recommendations. [13,14] An article's total score – quality - was calculated by the sum of

187  scores for each item, divided by the numbers of items and multiplied by 100%. [13,14] Overall,

188  the quality summary of appraised articles range from (0%-30%) Poor, (31%-50%) Fair, (51%-

189  70%) Good, (71%-90%) Very Good, and (>90%) Excellent. [13,14]

190

191  *Synthesis of Results*

192  A qualitative synthesis was conducted to report findings on test-retest reliability statistics. A meta-

193  analysis of Pearson's and Spearman's correlation was performed in Comprehensive Meta-

194  Analysis 3.3 software (Englewood, NJ). The meta-analyses were conducted using a random effect

8

195  model and the correlation coefficients were converted to z values. Heterogeneity was deemed

196  substantial if $I^2$ values were more than 50%. [19] A Meta-regression was planned to explore the

197  sources of unexplained heterogeneity by considering the following factors: a. neck pain with or

198  without radicular symptoms, b. acute or chronic, c. age and d. sex. Forest plots were created using

199  means and 95% confidence intervals for correlation coefficients. We summarize the main results

200  of the included articles based on the neck disorders, reported psychometric estimate and the study

201  quality ratings.

202

203  **RESULTS**

204  *Study Selection*

205  Our search yielded 123 articles. After removal of duplicates, 106 studies remained and were

206  screened using their title and abstract; leaving 28 articles selected for full-text review. Of these, 17

207  studies were considered eligible. [20,21,30–35,22–29] The flow of the study selection process is

208  presented in **Figure 1.**

209

210  *Study Characteristics*

211  The 16 eligible studies were conducted between 2006 and 2017 and included 1533 participants

212  with neck pain/disorders (mean of 96 participants per study). [20,21,30,32–35,22–29] Study size

213  ranged from 29 to 200 participants. A summary description of all the studies included is displayed

214  in **Table 1.** Concurrent validity was evaluated in 14 studies by comparing the difference of pain

215  intensity, disability and function scores with the score of GROC scales. Two studies [24,29]

216  examined the test-retest reliability of a 7-point and an 11-point GPE scale for patients with

217  whiplash-associated disorders (WAD). One study [22] examined whether occurrences of within-

9

218  and between-session changes were significantly associated with functional outcomes, pain, and

219  self-report of recovery in patients at discharge who were treated with manual therapy for

220  mechanical neck pain.

221

222  *COSMIN Risk of Bias rating and Quality appraisal of the Included Studies*

223  Regarding the risk of bias, all studies were rated as very good (**Table 2**). The quality of the studies

224  ranged from 88% to 96% (**Table 3**). The most common flaws were 1) lack of/inadequate sample

225  size calculations, 2) missing data (i.e. inadequate follow up), and 3) inconsistencies between the

226  data presented and hypothesis stated.

227

228  *Reported GROC scales*

229  The most commonly reported GROC scale (n=6 studies) was a 15-point scale with the most

230  frequent anchors being "-7 (a very great deal worse) to zero (about the same) to +7 (a very great

231  deal better)". A 7-point scale was reported in 5 studies, 11- and 5-point scales were reported in 2

232  studies and a 9-point scale in one study. The anchors in those scales varied greatly and are

233  presented in Table 1. Only 6 studies [24,29–31,33,34] reported full detail regarding the specific

234  questions asked of the patients with neck disorder when a GROC scale was administered. Those

235  questions that were reported are presented in **Box 1.**

236

237  *Reliability Measures*

238  Two studies were included that examined test-retest reliability of GPE for patients with WAD.

239  Kamper et al. (2010) [24]  examined the [time interval] test-retest reliability of an 11-point GPE

240  scale in 134 patients with chronic WAD and reported an Intra-class Correlation Coefficient (ICC)

10

241 of 0.99 (95% CI 0.99 to 0.99) at baseline, 0.96 (0.95 to 0.97) at 6 weeks, and 0.92 (0.89 to 0.94)

242 at 12 months. (**Table 4**). Ngo et al. (2010) assessed the test-retest reliability of a 7-point scale of

243 GPE in patients with acute WAD at 3 to 5 days. [29] The ICC and 95% confidence intervals (CI)

244 were used to determine the test–retest reliability of the two versions of the perceived recovery

245 questions using their original seven-item responses. Ngo et al. also computed weighted kappa

246 coefficients and 95% CI using quadratic weights to determine whether the distribution of responses

247 influenced the reliability as measured by the ICC. An ICC for general recovery of 0.70 (0.60 to

248 0.80) () and an ICC for neck pain questions of 0.80 (0.72 to 0.87) were found. A weighted Kappa

249 was also calculated (Kappa = 0.70 (0.42 to 0.98)) at six weeks for general recovery and at six

250 weeks Kappa = 0.80 (0.51 to 1.0) for neck pain questions (**Table 4**).

251

*Validity Measures*

253 We found 14 studies that examined concurrent validity measures between GROC and another PRO

254 (**Table 5**). Bjorklund et al. compared the validity of GROC with ProFitMAP-neck change scores

255 (moderate correlations: rho = 0.47, (p<0.05) and the Neck Disability Index (NDI) (moderate

256 correlations: rho = 0.59, (p<0.05) in patients with non-specific neck-shoulder pain.[30] Cleland et

257 al. compared the validity of GROC with NDI change scores (very weak correlations: $r = 0.19$) and

258 with Patient Specific Functional Scale change scores (PSFS) (very strong correlations: $r = 0.82$)

259 in 38 patients with cervical radiculopathy.[20] Cleland et al. compared the GROC with NDI

260 change scores (moderate correlations: $r = 0.58$) and with Numeric Pain Rating Scale (NPRS)

261 scores (moderate correlations: $r = 0.57$) in 137 patients with neck pain.[21] Farooq et al. compared

262 the GROC with the Urdu version of NDI change scores, and indicated moderate correlations $r =$

263 0.50 in 106 patients with neck pain.[36] Guzy et al. compared the GROC with NDI change scores

11

264 and reported moderate to strong correlations $r$ = -0.73 at two weeks and -0.56 at four weeks, in 95

265 patients with neck pain.[23] Jorritsma et al. compared the validity of GPE with Neck Pain and

266 Disability Scale change scores (NPAD) (moderate correlations: r = 0.49 (95% CI 0.30 to 0.64) in

267 patients with chronic non-specific neck pain. [32] Monticone et al. compared the GPE with

268 NeckPix change scores (strong correlations: rho = 0.69 to 0.82) in patients with chronic neck

269 pain.[33] Monticone et al. compared the GPE with the Italian version NDI change scores

270 (moderate correlations: Spearman's coefficient = 0.59) in patients with chronic neck pain. [34]

271 Shaheen et al. compared the validity of GROC with the Arabic version of NDI change scores and

272 indicated very strong correlations: $r$ coefficient = 0.81, in 70 patients with neck pain lasting more

273 than three months.[25] Takeshita et al. compared the validity of PGIC with the original NDI and

274 the Japanese version of NDI-J change scores and reported moderate correlations: $r$ coefficient =

275 0.47, and r = 0.59 in 130 patients with neck pain, cervical radiculopathy and/or cervical

276 myelopathy respectively.[26] Trouli et al. compared the validity of the GROC with the Greek

277 version of NDI change scores and reported weak correlations: $r$ coefficient = 0.30, in 68 patients

278 with neck pain.[27] Tuttle et al. compared the validity of GPE with NDI ($r$ coefficient range: 0.01

279 to 0.17; very weak correlations), with PSFS ($r$ coefficient range: 0.03 to 0.06; very weak

280 correlations), with pain intensity ($r$ coefficient range: 0.00 to 0.05; very weak correlations), and

281 with ROM ($r$ coefficient range: 0.00 to 0.03; very weak correlations), in 29 patients with neck pain

282 for more than two weeks.[28] Young et al. compared the validity of GROC with NDI change

283 scores and reported moderate correlations ($r$ coefficient = 0.52) in patients with mechanical neck

284 pain.

285

286 *Meta-Analysis and Meta-Regression of Correlations between Disability change scores and GROC*

287 *scores*

288 Five studies [21,23,32,35,36] of very good-to-excellent quality reported the Pearson correlation

289 coefficients between neck disability change scores and the GROC scores and were pooled together.

290 We found that GROC was positively correlated with disability change scores (r = 0.53, 95% CI:

291 0.47 to 0.59, $I^2$ = 0%). Six studies [25–28,30,34] of very good-to-excellent quality reported the

292 Spearman correlation coefficients between neck disability changes scores and the GROC scores

293 and were pooled together. We found that GROC was moderately correlated with disability change

294 scores (rho = 0.56, 95% CI: 0.41 to 0.68, $I^2$= 85%). The forest plots with correlation coefficients

295 with 95% CIs are presented in Figure 2-3. Our meta-regression showed that age was found as a

296 significant factor in influencing Fisher's Z scores (β = -0.034, 95% CI -0.05 to -0.01, p = 0.001).

297 The model explained 68% of the variance ($R^2 = 0.68$) (Figure 4).

298

299 *Area under the curve (AUC) – Sensitivity and Specificity*

300 Cook et al. [22] found that between-session NPRS- pain changes were associated with greater than

301 3-point change on the GROC at 96-hours (AUC=0.76). The pain change associated with GROC

302 was more specific (Specificity=79.2%, range: 62.2 - 91.1) than sensitive (Sensitivity=65.6%,

303 range: 57.9 to 74.6). Those with a 36.7% between-sessions change in pain were also 7.3 times

304 more likely to report an improvement of greater than 3 points change on the GROC than those

305 who did not achieve a 36.7% change in pain (**Table 4**).

306

307 **DISCUSSION**

13

308    This review has synthesized the current research from 17 studies that aimed to evaluate the

309    psychometric properties of GROC scales for patients with neck disorders, with the goal to provide

310    evidence for clinicians and researchers concerning its use within clinical practice and research.

311    From the 17 included studies, only 2 studies [24,29] reported test-retest reliability statistics of the

312    7- and 11-points item GPE scales for patients with WAD only. We were able to pool data from 12

313    studies regarding concurrent validity of GROC scales and neck disability change scores at one

314    time point after the interventions.[3] Themes influencing interpretation of the GROC were explored

315    in a study [31] that evaluated the factors that contribute to how patients respond to a question on

316    global perceived effect. This study found that treatment process, biomechanical performance, self-

317    efficacy and the nature of the condition may influence the responses on global perceived effect,

318    which is consistent with what we would expect for patients with neck pain. This suggests that

319    change is a complex multifactorial global concept. A strength of GROC is that it is intended as a

320    global assessment, and it can be assumed that it reflects the aspects of change important to the

321    individual patient.

322    Reliability can be defined as the degree to which a measure produces consecutive results

323    with the least amount of random error when the status of the population remains unchanged. The

324    reliability of GPE displayed an excellent test-retest reliability of ICC>0.90 over an interval of 6

325    weeks and 12 months for patients with WAD. Conducting an assessment with a long test-retest

326    interval (e.g. 12 months), can provide challenges as there is higher risk of individuals with WAD

327    being symptomatically unstable.[9] Determining if patients are symptomatically-stable can be

328    achieved by administering another PRO such as the Single Assessment Numeric Evaluation

329    (SANE)[37], however, the 7- and 11- points GPE scales still demonstrated good stability properties

14

330  at long test intervals (i.e., of 6 weeks and 12 months). Therefore, the measurements of the

331  reliability parameters of the GPE may be very useful during longer test intervals in clinical trials.

332  The psychometric property of validity is defined as the degree to which a PRO measures

333  what it is intended to measure. Pooled data from 11 studies overall suggest that post-treatment

334  changes of on validated disability outcome measures were moderately (Pearson's r = 0.51, 95%

335  CI: 0.43 to 0.58; Spearman's rho = 0.56, 95% CI: 0.41 to 0.68) correlated to change in perceived

336  effect) (Figure 2-3). This finding suggests that GROC scores taken at one point in time were related

337  to scores in pain and disability in patients with neck disorders, as measured by standardized

338  measures taken at 2 points in time. We identified one study [22] that found a 36.7% change in pain

339  for within- and between- session changes was associated with a 50% reduction in the NDI and an

340  improvement of >3 points on a 15-points GROC scale for patients with neck pain. This quantified

341  predictive change value may have clinical utility for use in clinical practice.

342  Previous studies [9,38] have indicated serious concerns about the conceptual validity of the

343  global rating of change. The review by Kamper et al.[9]  clearly showed that GROC was related

344  to final status more than change and was least related to baseline health status. This result

345  undermines the premise of what the global rating of change actually measures. For this reason, we

346  conclude that the 0.50 pooled correlation across 12 studies between the GROC and other PROM

347  change scores (e.g. NDI scores) may reflect a relationship between follow-up status and change

348  rather than supporting the contention that GROC actually measures change. This would also

349  explain why only 25% of the variation in GROC change scores was explained by changes scores

350  from a PROM change score measured at 2 points in time. In all studies, participants completed the

351  GROC scale at one time point after the intervention, and hence recall bias is a cause for concern.

352  However, another potential factor for moderate correlations is that the PROM, used as a

15

353 comparator, may not reflect the issues or priorities that are important to patients. Since no studies

354 compared a retrospective global assessment of the GROC to pre-post single item global PROM

355 e.g. the SANE, we do not know the extent to which these two factors contributed to moderate

356 correlation.

357 A unique aspect of this study was that it focused on global rating of change scales in a neck

358 pain patient population. Our study appraisal suggests that future studies concerning GROC should

359 include adequate sample sizes, maintain a rigorous follow up and report appropriate statistical error

360 estimates, since these were often inadequate. Various critical appraisal tools exist, and the

361 perspectives and ratings may differ across instruments. We used 2 different critical appraisal tools

362 to evaluate quality from 2 perspectives. The COSMIN risk of bias assessments reflects the level

363 of confidence in the conclusions and pooled estimates. The quality appraisal tool focuses on design

364 issues in the studies and reflects gaps in research designs that should be considered in interpretation

365 of current research and improved in future studies. Substantial heterogeneity was detected

366 ($I^2$>50%) in pooled Spearman's correlation coefficients which is a concern when pooling data. Our

367 univariate meta-regression analysis indicated that age across the studies explained 68% of the

368 variance (**Figure 4**). Other factors such as type of neck pain (with or without radicular symptoms),

369 acute or chronic and sex did not explain the remaining heterogeneity (not statically significant).

370 Furthermore, the scope of our literature search was focused on identifying full-text papers written

371 in English only.

372 While this study included 16 studies, only 2 of these reported reliability statistics for GROC

373 scales for patients with chronic WAD. Therefore, the applicability of our study is mostly limited

374 to patients with chronic WAD. For validity measurements, GROC scales were mostly investigated

375 by correlation analyses to evaluate the external responsiveness of another PRO measure over a

16

376 specific time point. From our meta-analysis, we can be confident that the GROC scores were

377 moderately correlated with neck disability change scores. However, more robust psychometric

378 design studies to test the measurement properties of GROC scales as the primary outcome of

379 investigation are highly needed. Future studies should aim to test to what extent the different range

380 of items (e.g. 7-point scale vs 11-point scale), the anchors (e.g. much worse vs much better) may

381 affect the measurement properties of GROC scales for patients with neck disorders.

382

**CONCLUSIONS**

384 This study found excellent quality evidence of very good to excellent test-retest reliability of GPE

385 for patients with WAD. Evidence of very good to excellent quality studies found that GROC scores

386 are moderately correlated to an external criterion PROM measure measured pre-post treatment in

387 patients with neck disorders. Studies addressing the optimal form of GROC scales for patients with

388 neck disorders or comparing the GROC to other options for single-item global assessment of

389 change were not found.

390

**Authors' contributions**

392 PB contributed significantly to conception and design of the study, data extraction, critical

393 appraisal, interpretation of data and drafting of the manuscript. GN, and RF were involved in

394 literature search, critical appraisal and interpretation of data and drafting. GN was involved in

395 critical appraisal and drafting. JM was also involved in the conception and design of the study,

396 drafting, and revised the manuscript for important intellectual content. JM and CATWAD were

397 involved in the drafting and review of the manuscript. All authors have given their final approval

398 on the manuscript to be published

17

399

**Declarations**

**Ethics approval and consent to participate**

Not applicable

**Consent for publication**

Not applicable

**Availability of data and material**

Data sharing is not applicable to this article as no datasets were generated or analyzed during the

current study

**Funding Statement**

**Competing Interest Statement**

None to report

**References**

1    Murray CJL, Abraham J, Ali MK, *et al.* The State of US health, 1990-2010: Burden of diseases,

injuries, and risk factors. *JAMA - J Am Med Assoc* Published Online First: 2013.

doi:10.1001/jama.2013.13805

2    Fejer R, Kyvik KO, Hartvigsen J. The prevalence of neck pain in the world population: A

systematic critical review of the literature. Eur. Spine J. 2006. doi:10.1007/s00586-004-0864-4

3    Hogg-Johnson S, van der Velde G, Carroll LJ, *et al.* The Burden and Determinants of Neck Pain in

the General Population. Results of the Bone and Joint Decade 2000-2010 Task Force on Neck

18

423    Pain and Its Associated Disorders. *J Manipulative Physiol Ther* Published Online First: 2009.

424    doi:10.1016/j.jmpt.2008.11.010

425    4    Bobos P, Nazari G, Palimeris S, *et al.* The contribution of health and psychological factors in

426    patients with chronic neck pain and disability: A cross-sectional study. *J Clin Diagnostic Res*

427    Published Online First: 2018. doi:10.7860/JCDR/2018/31284.11203

428    5    Macdermid JC, Walton DM, Bobos P, *et al.* The Open Orthopaedics Journal A Qualitative

429    Description of Chronic Neck Pain has Implications for Outcome Assessment and Classification.

430    *Open Orthop J* 2016;**10**:746–56. doi:10.2174/1874325001610010746

431    6    Treleaven J. Sensorimotor disturbances in neck disorders affecting postural stability, head and eye

432    movement control-Part 2: Case studies. *Man Ther* 2008;**13**:266–75.

433    doi:10.1016/j.math.2007.11.002

434    7    Cohen SP. Epidemiology, diagnosis, and treatment of neck pain. *Mayo Clin Proc* 2015;**90**:284–99.

435    doi:10.1016/j.mayocp.2014.09.008

436    8    Bobos P, MacDermid JC, Walton DM, *et al.* Patient-Reported Outcome Measures Used for Neck

437    Disorders: An Overview of Systematic Reviews. *J Orthop Sport Phys Ther* 2018;**48**:775–88.

438    doi:10.2519/jospt.2018.8131

439    9    Kamper SJ, Maher CG, Mackay G. Global Rating of Change Scales: A Review of Strengths and

440    Weaknesses and Considerations for Design. *J Man Manip Ther* 2009;**17**:163–70.

441    doi:10.1002/mus.21062

442    10    Nazari G, Bobos P, MacDermid JC, *et al.* The Effectiveness of Instrument-Assisted Soft Tissue

443    Mobilization in Athletes, Participants Without Extremity or Spinal Conditions, and Individuals

444    with Upper Extremity, Lower Extremity, and Spinal Conditions: A Systematic Review. *Arch Phys*

445    *Med Rehabil* Published Online First: February 2019. doi:10.1016/j.apmr.2019.01.017

446    11    Bobos P, Nazari G, Szekeres M, *et al.* The effectiveness of joint-protection programs on pain,

447    hand function, and grip strength levels in patients with hand arthritis: A systematic review and

448    meta-analysis. *J Hand Ther* 2018;**32**:194–211. doi:10.1016/j.jht.2018.09.012

19

449  12  Nazari G, Bobos P, MacDermid JC, *et al.* Psychometric properties of the Zephyr bioharness

450     device: A systematic review. *BMC Sports Sci Med Rehabil* 2018;**10**. doi:10.1186/s13102-018-

451     0094-4

452  13  Law MC, MacDermid J. *Evidence-based rehabilitation : a guide to practice.* Thorofare, NJ: :

453     Slack Incorporated 2014.

454  14  Roy JS, Desmeules F, MacDermid JC. Psychometric properties of presenteeism scales for

455     musculoskeletal disorders: A systematic review. *J Rehabil Med* Published Online First: 2011.

456     doi:10.2340/16501977-0643

457  15  Wuensch KL, Evans JD. Straightforward Statistics for the Behavioral Sciences. *J Am Stat Assoc*

458     Published Online First: 2006. doi:10.2307/2291607

459  16  Cohen J. Statistical power analysis for the behavioral sciences. Stat. Power Anal. Behav. Sci.

460     1988. doi:10.1234/12345678

461  17  Mokkink LB, de Vet HCW, Prinsen CAC, *et al.* COSMIN Risk of Bias checklist for systematic

462     reviews of Patient-Reported Outcome Measures. *Qual Life Res* Published Online First: 2018.

463     doi:10.1007/s11136-017-1765-4

464  18  Terwee CB, Mokkink LB, Knol DL, *et al.* Rating the methodological quality in systematic reviews

465     of studies on measurement properties : a scoring system for the COSMIN checklist. 2012;:651–7.

466     doi:10.1007/s11136-011-9960-1

467  19  Higgins JPT, Thompson SG, Deeks JJ, *et al.* Measuring inconsistency in meta-analyses. *BMJ*

468     Published Online First: 2003. doi:10.1136/bmj.327.7414.557

469  20  Cleland J, Fritz J, Whitman J, *et al.* The reliability and construct validity of the Neck Disability

470     Index and Patient Specific Functional Scale. *Spine (Phila Pa 1976)* 2006;**31**:598–602.

471  21  Cleland JA, Childs JD, Whitman JM. Psychometric Properties of the Neck Disability Index and

472     Numeric Pain Rating Scale in Patients With Mechanical Neck Pain. *Arch Phys Med Rehabil*

473     2008;**89**:69–74. doi:10.1016/j.apmr.2007.08.126

474  22  Cook C, Lawrence J, Michalak K, *et al.* Is there preliminary value to a within- and/or between-

20

475     session change for determining short-term outcomes of manual therapy on mechanical neck pain?

476     *J Man Manip Ther* 2014;**22**:173–80. doi:10.1179/2042618614y.0000000071

477   23   Guzy G, Vernon H, Polczyk R, *et al.* Psychometric validation of the authorized Polish version of

478     the Neck Disability Index. *Disabil Rehabil* 2013;**35**:2132–7. doi:10.3109/09638288.2013.771706

479   24   Kamper SJ, Ostelo RWJG, Knol DL, *et al.* Global Perceived Effect scales provided reliable

480     assessments of health transition in people with musculoskeletal disorders, but ratings are strongly

481     influenced by current status. *J Clin Epidemiol* 2010;**63**:760-766.e1.

482     doi:10.1016/j.jclinepi.2009.09.009

483   25   Shaheen AAM, Omar MTA, Vernon H. Cross-cultural adaptation, reliability, and validity of the

484     arabic version of neck disability index in patients with neck pain. *Spine (Phila Pa 1976)*

485     2013;**38**:609–15. doi:10.1097/BRS.0b013e31828b2d09

486   26   Takeshita K, Hosono N, Kawaguchi Y, *et al.* Validity, reliability and responsiveness of the

487     Japanese version of the Neck Disability Index. *J Orthop Sci* 2013;**18**:14–21. doi:10.1007/s00776-

488     012-0304-y

489   27   Trouli MN, Vernon HT, Kakavelakis KN, *et al.* Translation of the Neck Disability Index and

490     validation of the Greek version in a sample of neck pain patients. *BMC Musculoskelet Disord*

491     2008;**9**:1–8. doi:10.1186/1471-2474-9-106

492   28   Tuttle N, Laakso L, Barrett R. Change in impairments in the first two treatments predicts outcome

493     in impairments, but not in activity limitations, in subacute neck pain: An observational study. *Aust*

494     *J Physiother* 2006;**52**:281–5. doi:10.1016/S0004-9514(06)70008-3

495   29   Ngo Trung, Stupar Maja, Coˆteˊ Pierre, Boyle Eleanor, Shearer Heather. A study of the test –

496     retest reliability of the self-perceived general recovery and self-perceived change in neck pain

497     questions in patients with recent whiplash-associated disorders. 2010;:957–62.

498     doi:10.1007/s00586-010-1289-x

499   30   Björklund M, Wiitavaara B, Heiden M. Responsiveness and minimal important change for the

500     ProFitMap-neck questionnaire and the Neck Disability Index in women with neck–shoulder pain.

21

501   *Qual Life Res* 2017;**26**:161–70. doi:10.1007/s11136-016-1373-8

502   31   Evans R, Bronfort G, Maiers M, *et al.* '" I know it " s changed ''': a mixed-methods study of the

503   meaning of Global Perceived Effect in chronic neck pain patients. 2014;:888–97.

504   doi:10.1007/s00586-013-3149-y

505   32   Jorritsma W, Dijkstra PU, De Vries GE, *et al.* Detecting relevant changes and responsiveness of

506   Neck Pain and Disability Scale and Neck Disability Index. *Eur Spine J* 2012;**21**:2550–7.

507   doi:10.1007/s00586-012-2407-8

508   33   Monticone M, Frigau L, Vernon H, *et al.* Responsiveness and minimal important change of the

509   NeckPix© in subjects with chronic neck pain undergoing rehabilitation. *Eur Spine J*

510   2018;**27**:1324–31. doi:10.1007/s00586-017-5343-9

511   34   Monticone M, Ambrosini E, Vernon H, *et al.* Responsiveness and minimal important changes for

512   the Neck Disability Index and the Neck Pain Disability Scale in Italian subjects with chronic neck

513   pain. *Eur Spine J* 2015;**24**:2821–7. doi:10.1007/s00586-015-3785-5

514   35   Young BA, Walker MJ, Strunce JB, *et al.* Responsiveness of the Neck Disability Index in patients

515   with mechanical neck disorders. *Spine J* 2009;**9**:802–8. doi:10.1016/j.spinee.2009.06.002

516   36   Farooq MN, Mohseni-Bandpei MA, Gilani SA, *et al.* Urdu version of the neck disability index: A

517   reliability and validity study. *BMC Musculoskelet Disord* 2017;**18**:1–11. doi:10.1186/s12891-017-

518   1469-5

519   37   Williams GN, Gangel TJ, Arciero RA, *et al.* Comparison of the single assessment numeric

520   evaluation method and two shoulder rating scales. Outcomes measures after shoulder surgery. *Am*

521   *J Sports Med* 1999;**27**:214–21. doi:10.1177/03635465990270021701

522   38   Schmitt J, Abbott JH. Global Ratings of Change Do Not Accurately Reflect Functional Change

523   Over Time in Clinical Practice. *J Orthop Sport Phys Ther* 2015;**45**:106–11.

524   doi:10.2519/jospt.2015.5247

525

526

22

**Table 1**. Study Characteristics

| Study | Population | Setting | Sample Size | Properties Evaluated | GROC evaluated | Interval |
|-------|-----------|---------|-------------|----------------------|----------------|----------|
| Bjorklund et al (2017) | Women with non-specific neck-shoulder pain | Not specified | 104 | Validity (correlation) Between NDI and GRoC | GRoC 7-points 1. Very much worse; 2. Much worse; 3. Minimally worse; 4. No change; 5. Minimally improved; 6. Much improved; 7. Very much improved. | GRoC scale administered only after intervention at one time point (1 week) |
| Cleland et al (2006) | Patients with cervical radiculopathy | Hospital | 38 | Validity (correlation) Between NDI and GRoC Between PSFS and GRoC | GRoC 15-points -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at follow up. Within a week over the period of 7 weeks. |
| Cleland et al. (2008) | Patients with neck pain only | 5 Outpatient physical therapy clinics | 137 | Validity (correlation) Between NDI and GRoC Between NPRS and GRoC | GRoC 15-points -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at follow up. Within a week |
| Cook et al (2014) | Patients with any neck pain | Academic locations in Northeast Ohio | 56 | ROC curves and AUC to measure sensitivity and specificity. Binomial logistic regression analysis was also calculated to determine overall effect. | GRoC 15-points -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | Baseline and at follow up 48- and 96-hours post baseline |
| Farooq et al. (2017) | Patients with neck pain | Physical therapy clinics | 106 | Validity (correlation) Between NDI-U and GRoC | GRoC 15-points -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at three weeks after intervention |
| Guzy et al. (2013) | Patients with neck pain | Outpatient rehabilitation clinic | 95 | Validity (correlation) Between NDI-P and GRoC | GRoC 7-points 'complete recovery'' over ''no change'' to ''my complaints are worse than ever'' | GRoC scale was completed at 2 weeks and at 4 weeks |
| Jorritsma et al. (2012) | Patients with chronic non-specific neck pain | Tertiary university center for rehabilitation | 76 | Validity (correlation) Between NDI and GRoC Between NPAD and GRoC | GPE 7-points 3 (completely recovered) to zero (no change) to -3 (worse than ever) | After completion the program varying from 3 to 5 months patients filled the GPE |
| Kamper et al. (2010) | Patients with any whiplash-associated disorder. | Physical therapy clinics | 134 | Test-retest reliability | GPE 11-points -5 (vastly worse) to zero (unchanged) to +5 (completely recovered) | Baseline, 6 weeks and 12 months |
| Monticone et al. 2017 | Patients with chronic neck pain | Outpatient Rehabilitation Unit | 153 | Validity (correlation) Between NeckPix and GPE | GPE 5-points (helped a lot = 1, helped = 2), one no change level (helped only a little = 3), and two worsening levels (did not help = 4, made things worse = 5) | At the end of treatment (8 weeks) and one year before follow-up |
| Monticone et al. 2015 | Patients with chronic neck pain | Outpatient Rehabilitation Unit | 200 | Validity (correlation) Between NDI and GPE Between NPDS and GPE | GPE 5-points (helped a lot = 1, helped = 2), one no change level (helped only a little = 3), and two worsening levels (did not help = 4, made things worse = 5) | At the end of treatment 8 weeks |

527

528

23

| | | | | | | |
|---|---|---|---|---|---|---|
| Ngo et al. (2010) | Patients with WAD. Most participants (69.6%) had grade II WAD. | Interviewed by person or by telephone in Ontario | 46 | Test-retest reliability | GPE 7-points 1. General recovery question Completely better Much improved Slightly improved No change Slightly worse Much worse Worse than ever 2. Change in neck pain question: very much better, better, slightly better, no change, slightly worse, worse, or very much worse | 3-5 days |
| Shaheen et al. (2015) | Patients with neck pain lasting more than 3 months | 3 primary health centers | 70 | Validity (correlation) Between NDI-Ar and GRoC | GRoC 15-points -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | 1 week |
| Takeshita et al. (2014) | Patients with neck pain, cervical radiculopathy and/or cervical myelopathy | Variety of clinics and hospital settings | 130 | Validity (correlation) Between NDI-J and GRoC | PGIC 7-points much better, better, slightly better, unchanged, slightly worse, worse and much worse | Over 8 weeks |
| Trouli et al. (2008) | Patients with neck pain | Primary healthcare clinic | 68 | Validity (correlation) Between NDI-Gr and GRoC | GRoC 15-points -7 (a very great deal worse) to -1 (almost the same, hardly any worse at all) and from 7 (a very great deal better) to 1 (almost the same, hardly any better at all) | Within 2 months but 1 week for test-retest |
| Tuttle et al. (2006) | Patients with neck pain for more than 2 weeks | Private physiotherapy clinics | 29 | Validity (correlation) Between NDI and GPE Between PSFS and GPE Between VAS and GPE Between ROM and GPE | GPE 11-points −5 is vastly worse and +5 is completely recovered | 6 weeks |
| Young et al. (2009) | Patients presenting with mechanical neck pain | Outpatient physical therapy clinics. | 91 | Validity (correlation) | GRoC 15-points -7 (''a very great deal worse'') to 0 (''about the same'') to +7 (''a very great deal better'') | 3 weeks |

529

**TABLE 3.** *Quality Appraisal for Clinical Measurement Research Reports Evaluation Form*

| Study | Item Evaluation Criteria* | | | | | | | | | | | | Total (%) | Quality Summary |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | | |
| Bjorklund et al (2017) | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Cleland et al. (2008) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Trouli et al. (2008) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Tuttle et al. (2006) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Kamper et al. (2010) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Cook et al (2014) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 92 | Excellent |

531 **TABLE 2.** Summary of Psychometric Properties Reported in Studies and COSMIN Risk of Bias (RoB)
532 and Quality studies

| Study | Psychometric Properties Reported | COSMIN RoB | COSMIN Rating*§ (Criteria) | Quality of Studies** (QACMRR) |
|---|---|---|---|---|
| Bjorklund et al (2017) | Validity (correlation) | Very Good | ? | Excellent |
| Cleland et al (2006) | Validity (correlation) | Very Good | + | Excellent |
| Cleland et al. (2008) | Validity (correlation) | Very Good | - | Excellent |
| Cook et al (2014) | Sensitivity Specificity | Very Good Very Good | + | Excellent |
| Farooq et al. (2017) | Validity (correlation) | Very Good | + | Excellent |
| Guzy et al. (2013) | Validity (correlation) | Very Good | ? | Very good |
| Jorritsma et al. (2012) | Validity (correlation) | Very Good | ? | Excellent |
| Kamper et al. (2010) | Test-retest reliability | Very Good | + | Excellent |
| Monticone et al. (2017) | Validity (correlation) | Very Good | ? | Excellent |
| Monticone et al. (2015) | Validity (correlation | Very Good | ? | Excellent |
| Ngo et al. (2010) | Test-retest reliability | Very Good | + | Excellent |
| Shaheen et al. (2015) | Validity (correlation) | Very Good | ? | Excellent |
| Takeshita et al. (2014) | Validity (correlation) | Very Good | ? | Very good |
| Trouli et al. (2008) | Validity (correlation) | Very Good | + | Excellent |
| Tuttle et al. (2006) | Validity (correlation) | Very Good | ? | Excellent |
| Young et al. (2009) | Validity (correlation) | Very Good | ? | Excellent |

533 COSMIN, Consensus-based Standards for the Selection of health Measurement Instruments, Criteria for good measurement
534 properties: '+' sufficient; '-'insufficient; '?' indeterminate. §§ The grading for the quality of the evidence based on the modified
535 GRADE approach is not applicable. **Quality Appraisal for Clinical Measurement Research Reports Evaluation Form
536 (QACMRR).

25

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Jorritsma et al. (2012) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Cleland et al (2006) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Monticone et al. (2017) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Monticone et al. (2015) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Ngo et al. (2010) | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 1 | 2 | 92 | Excellent |
| Shaheen et al. (2013) | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 92 | Excellent |
| Farooq et al. (2017) | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 92 | Excellent |
| Young et al. (2009) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Guzy et al. (2013) | 2 | 2 | 1 | 2 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 88 | Very good |
| Takeshita et al. (2014) | 2 | 2 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 88 | Very good |

537

538 *Item Evaluation Criteria: 1. Thorough literature review to define the research question; 2. Specific inclusion/exclusion

539 criteria; 3. Specific hypotheses; 4. Appropriate scope of psychometric properties; 5. Sample size; 6. Follow-up; 7. The

540 authors referenced specific procedures for administration, scoring, and interpretation of procedures; 8. Measurement

541 techniques were standardized; 9. Data were presented for each hypothesis; 10. Appropriate statistics-point estimates; 11.

542 Appropriate statistical error estimates; 12. Valid conclusions and clinical recommendations.

543 Total score = (sum of subtotals ÷ 24 × 100). If for a specific paper an item is deemed NA (Not Applicable), then, Total score

544 = (sum of subtotals ÷ (2 × number of Applicable items) × 100).

545 NA – Not Applicable. The subsections no. 6, asks for percentage of retention/follow up. This subsection only applies to

546 reliability test-retest studies

547 Quality Summary: Poor (0%-30%), Fair (31%-50%), Good (51%-70%), Very good (71%-90%), Excellent (>90%):

548

**TABLE 5.** SUMMARY OF VALIDITY PROPERTIES OF GRoC SCALES

| Study | Type of Validity | Validity Estimates | COSMIN | Quality |
|---|---|---|---|---|
| | Spearman's correlation between the change scores | | Very Good | Excellent |

Bjorklund et al (2017)

**TABLE 4.** SUMMARY OF RELIABILITY PROPERTIES OF GROC SCALES

GRoC and ProFitMap-neck — rho = 0.47, (p<0.05)
GRoC and NDI — rho = 0.59, (p<0.05)

| Study | Type of Reliability | Reliability Estimates | COSMIN | Quality |
|---|---|---|---|---|
| Cleland et al. (2006) | Test-retest; Correlations (Pearson r) between change scores; NDI and GRoC; PSFS and GRoC | Intra-class correlation coefficients (ICC) 0.99 (0.99–0.99) – baseline 0.96 (0.95 – 0.97) – at six weeks 0.92 (0.89 – 0.94) at twelve months. r = 0.58 r = 0.57 | Very Good | Very Good |
| Cleland et al. (2008) | Correlations (Pearson r) between change scores; NDI and GRoC; NRS and GRoC | Intra-class correlation coefficients (ICC) 0.70 (0.60–0.80) – at six weeks (General recovery) 0.80 (0.72–0.87) – at six weeks (neck pain questions) | Very Good | Excellent |
| Cook et al. (2014) | Receiver operator characteristics (ROC); Within-session change; Between-session change; Between session change of Pain and GROC; Sensitivity; Specificity | Weighted Kappa AUC = 0.61 AUC = 0.76, >36.7% change in pain 0.70 (0.42–0.98) – at six weeks (General recovery) 0.80 (0.51–1.0) – at six weeks (neck pain questions) Odds ratio = 7.3 (2.1, 24.7) 65.6% (57.9–74.6) 79.2% (62.2–91.1) 0.85 (0.64–1) when ''recovered'' was defined ''completely better'' 0.81 (0.64–0.99) when defined as ''completely better'' or ''much improved'' | Very Good | Excellent |
| Farooq et al. (2017) | Correlations (Pearson r); NDI-U | r = 0.50 | Very Good | Excellent |
| Guzy et al. (2013) | Correlations (Pearson r); NDI vs GROC | Dichotomized response options for change in neck pain questions (K statistics) Two-week interval (r = 0.73) Four-week interval (r = 0.56) | Very Good | Very good |
| Jorritsma et al. (2012) | Test-retest; Correlation between change scores of NPAD and GPE | 0.46 (0.20–0.74) when ''recovered'' was defined as ''very much better'' r = 0.49 (95 % CI 0.30–0.64) 0.80 (0.62–0.99) when defined as ''very much better'' or ''better'' | Very Good | Excellent |
| Monticone et al. (2017) | Correlations (Spearman) between change scores of the NeckPix© and GPE | Recall questions (K statistics) the kappa coefficient was 1 for participants who remembered their previous answers to the general recovery question; 0.88 (0.64–1) for those who did not remember and 0.50 (0.02–0.98) for participants who were not asked the question. rho = 0.69–0.82 | Very Good | Excellent |
| Monticone et al. (2015) | Correlation (Spearman) between change scores; NDI-I and GPE; NDPS and GPE | rho = 0.71, p<0.01 rho = 0.59, p<0.01 The kappa coefficient was 1 for participants who remembered their previous answers to the change in neck pain question; 0.74 (0.41–1) for those who did not remember and 0.66 (0.22–1) for participants who were not asked the question. | Very Good | Excellent |
| Shaheen et al. (2013) | Correlations (Spearman's); NDI-Ar and GROC | rho = 0.81, p<0.01 | Very Good | Excellent |
| Takeshita et al. (2014) | Correlations; NDI and PGIC; NDI-J and PGIC | Spearman (rho) rho = 0.47, p<o.oo1 rho = 0.59, p<o.oo1 | Very Good | Very good |
| Trouli et al. (2008) | Correlation (Spearman's); GROC vs Gr-NDI | rho = 0.30, p=0.02 | Very Good | Excellent |

549

27

| | | | Very Good | Excellent |
|---|---|---|---|---|
| | Correlations (Spearman's) | | | |
| | NDI vs GPE (post 1, minus pre-1) | | | |
| | NDI vs GPE (post 2, minus pre-1) | | | |
| | NDI vs GPE (post 2, minus pre-2) | rho = 0.17 | | |
| | | rho = 0.01 | | |
| | PSFS vs GPE (post 1, minus pre-1) | rho = 0.03 | | |
| | PSFS vs GPE (post 2, minus pre-1) | rho = 0.06 | | |
| | PSFS vs GPE (post 2, minus pre-2) | rho = 0.03 | | |
| Tuttle et al. (2006) | | rho = 0.03 | | |
| | Pain Intensity (post 1, minus pre-1) | rho = 0.00 | | |
| | Pain Intensity (post 2, minus pre-1) | rho = 0.05 | | |
| | Pain Intensity (post 2, minus pre-2) | rho = 0.01 | | |
| | | rho = 0.03 | | |
| | Total ROM (post 1, minus pre-1) | rho = 0.01 | | |
| | Total ROM (post 2, minus pre-1) | rho = 0.00 | | |
| | Total ROM (post 2, minus pre-2) | | | |
| Young et al. (2009) | Correlations (Pearson's) between change scores NDI and GRoC | r =0.52 (p<0.01) | Very Good | Excellent |

550

551

552

553

554

555

556

557

558

559

560

561

562

563 **Box 1.** Questions of Global Rating of Change (GROC) scales

| Author | GROC item- scale | Patients with neck disorders were asked: |
|---|---|---|
| Bjorklund et al. (2017) | GROC 7-points | *"Compared to before the treatment of the study started, my overall status is now"* <br> *"Compared to before the treatment of the study started, my status regarding my neck–shoulder problem is now"* |
| Evans et al (2014) | GPE 9-points | *"Overall, how much has your neck pain changed since you started treatment in the study?"* |
| Kamper et al. (2010) | GPE 11-points | *"With respect to your whiplash injury how would you describe yourself now compared to immediately after your accident"* |
| Monticone et al. (2017) | GPE 5-points | *"Overall, how much did the treatment you received help your fear of movement due to current neck pain?* <br> *"Overall, how much did the treatment you delivered help your subject's fear of movement due to her/ his current neck pain?"* |
| Monticone et al. (2015) | GPE 5-points | *"Overall, how much did the treatment you received help your neck problem?"* |
| Ngo et al. (2010) | GPE 7-points | *"How well do you feel you are recovering from your injuries?"* <br> *"How do you feel your neck pain has changed since the injury?"* |

564

565

566 **Figure 1.** Flow diagram of included studies

567 **Figure 2**. Meta-analysis of Pearson's correlation coefficients between neck disability change scores and
568 GROC scores in patients with neck disorders based on 5 very good to excellent quality studies.

569 **Figure 3**. Meta-analysis of Spearman's correlation coefficients between neck disability change scores
570 and GROC scores in patients with neck disorders based on 6 very good to excellent quality studies.

571 **Figure 4**. Random effects univariate meta-regression between age and the Fisher's Z estimates. Each circle
572 represents a study and the size of the circle indicates the influence of that study on the model. The
573 regression prediction is illustrated by the straight line and the curved lines represent the 95% confidence
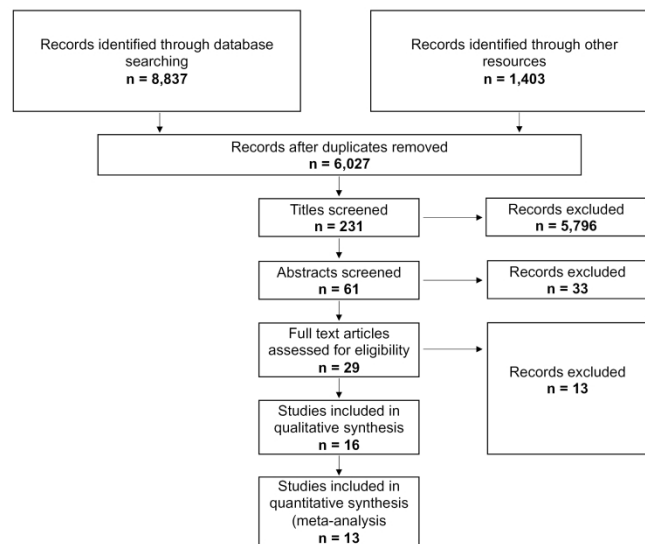574 intervals. Age explained 68% of the variance in the model ($R^2=0.68$).

29

Figure 1. Flow diagram of included studies
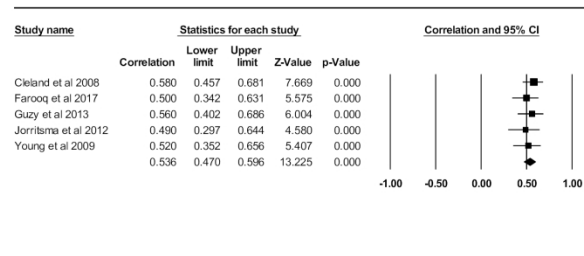
338x190mm (300 x 300 DPI)

Figure 2. Meta-analysis of Pearson's correlation coefficients between neck disability change scores and GROC scores in patients with neck disorders based on 5 very good to excellent quality studies.

215x279mm (300 x 300 DPI)

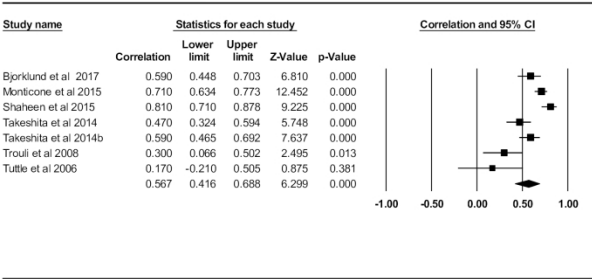| Study name | Statistics for each study | | | | | Correlation and 95% CI |
|---|---|---|---|---|---|---|
| | Correlation | Lower limit | Upper limit | Z-Value | p-Value | |
| Bjorklund et al 2017 | 0.590 | 0.448 | 0.703 | 6.810 | 0.000 | |
| Monticone et al 2015 | 0.710 | 0.634 | 0.773 | 12.452 | 0.000 | |
| Shaheen et al 2015 | 0.810 | 0.710 | 0.878 | 9.225 | 0.000 | |
| Takeshita et al 2014 | 0.470 | 0.324 | 0.594 | 5.748 | 0.000 | |
| Takeshita et al 2014b | 0.590 | 0.465 | 0.692 | 7.637 | 0.000 | |
| Trouli et al 2008 | 0.300 | 0.066 | 0.502 | 2.495 | 0.013 | |
| Tuttle et al 2006 | 0.170 | -0.210 | 0.505 | 0.875 | 0.381 | |
| | 0.567 | 0.416 | 0.688 | 6.299 | 0.000 | |

Figure 3. Meta-analysis of Spearman's correlation coefficients between neck disability change scores and GROC scores in patients with neck disorders based on 6 very good to excellent quality studies.

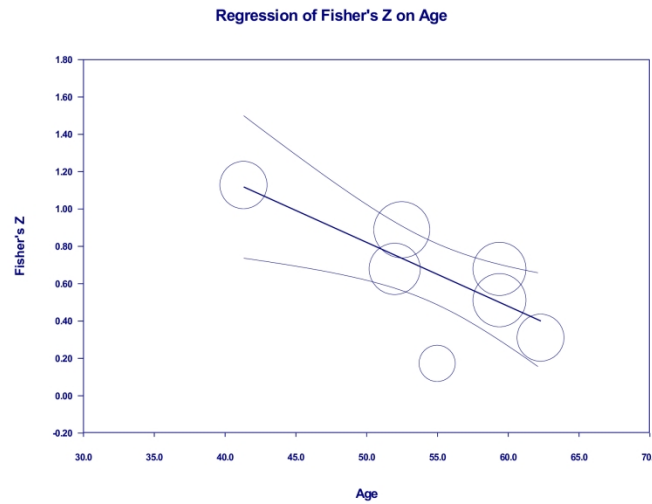215x279mm (300 x 300 DPI)

**Regression of Fisher's Z on Age**

Figure 4. Random effects univariate meta-regression between age and the Fisher's Z estimates. Each circle represents a study and the size of the circle indicates the influence of that study on the model. The regression prediction is illustrated by the straight line and the curved lines represent the 95% confidence intervals. Age explained 68% of the variance in the model (R2=0.68).

215x279mm (300 x 300 DPI)

**Appendix 1: Search terms**

MEDLINE-OVID
1. exp "outcome and process assessment (health care)"/ or "outcome assessment (health care)"/ or treatment outcome/
2. outcome?.ti.
3. exp "Range of Motion, Articular"/
4. Pain Measurement/
5. exp disability evaluation/
6. "Recovery of Function"/
7. Questionnaires/
8. self-report.tw.
9. ((impairment or disability or function) adj2 (measure? or scale? or evaluation?)).tw.
10. range of motion.tw.
11. (strength adj2 (measure? or scale? or evaluation?)).tw.
12. (outcome? adj2 (measure* or scale? or indicator?)).tw.
13. or/1-12
14. "reproducibility of results"/
15. exp "Sensitivity and Specificity"/
16. reliability.mp.
17. validity.mp.
18. responsiveness.mp.
19. Psychometrics/
20. rasch.mp.
21. factor analysis, statistical/
22. factor analysis.tw.
23. differential functioning.mp.
24. (validity or validation).mp. [mp=title, original title, abstract, name of substance word, subject heading word, unique identifier]
25. (validity or validation).mp.
26. item difficulty.mp.
27. translation.tw.
28. or/14-27
29. 13 and 28
30. Neck Pain/
31. exp Brachial Plexus Neuropathies/
32. exp neck injuries/ or exp whiplash injuries/
33. cervical pain.mp.
34. neckache.mp.
35. whiplash.mp.
36. cervicodynia.mp.
37. cervicalgia.mp.
38. brachialgia.mp.
39. brachial neuritis.mp.
40. brachial neuralgia.mp.
41. neck pain.mp.

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

42. neck injur*.mp.
43. brachial plexus neuropath*.mp.
44. brachial plexus neuritis.mp.
45. thoracic outlet syndrome/ or cervical rib syndrome/
46. Torticollis/
47. exp brachial plexus neuropathies/ or exp brachial plexus neuritis/
48. cervico brachial neuralgia.ti,ab.
49. cervicobrachial neuralgia.ti,ab.
50. (monoradicul* or monoradicl*).tw.
51. or/30-50
52. exp headache/ and cervic*.tw.
53. exp genital diseases, female/
54. genital disease*.mp.
55. or/53-54
56. 52 not 55
57. 51 or 56
58. neck/
59. neck muscles/
60. exp cervical plexus/
61. exp cervical vertebrae/
62. atlanto-axial joint/
63. atlanto-occipital joint/
64. Cervical Atlas/
65. spinal nerve roots/
66. exp brachial plexus/
67. (odontoid* or cervical or occip* or atlant*).tw.
68. axis/ or odontoid process/
69. Thoracic Vertebrae/
70. cervical vertebrae.mp.
71. cervical plexus.mp.
72. cervical spine.mp.
73. (neck adj3 muscles).mp.
74. (brachial adj3 plexus).mp.
75. (thoracic adj3 vertebrae).mp.
76. neck.mp.
77. (thoracic adj3 spine).mp.
78. (thoracic adj3 outlet).mp.
79. trapezius.mp.
80. cervical.mp.
81. cervico*.mp.
82. 80 or 81
83. exp genital diseases, female/
84. genital disease*.mp.
85. exp *Uterus/
86. 83 or 84 or 85
87. 82 not 86

88. 58 or 59 or 60 or 61 or 62 or 63 or 64 or 65 or 66 or 67 or 68 or 69 or 70 or 71 or 72 or 73 or
74 or 75 or 76 or 77 or 78 or 79 or 87
89. exp pain/
90. exp injuries/
91. pain.mp.
92. ache.mp.
93. sore.mp.
94. stiff.mp.
95. discomfort.mp.
96. injur*.mp.
97. neuropath*.mp.
98. or/89-97
99. 88 and 98
100. Radiculopathy/
101. exp temporomandibular joint disorders/ or exp temporomandibular joint dysfunction
syndrome/
102. myofascial pain syndromes/
103. exp "Sprains and Strains"/
104. exp Spinal Osteophytosis/
105. exp Neuritis/
106. Polyradiculopathy/
107. exp Arthritis/
108. Fibromyalgia/
109. spondylitis/ or discitis/
110. spondylosis/ or spondylolysis/ or spondylolisthesis/
111. radiculopathy.mp.
112. radiculitis.mp.
113. temporomandibular.mp.
114. myofascial pain syndrome*.mp.
115. thoracic outlet syndrome*.mp.
116. spinal osteophytosis.mp.
117. neuritis.mp.
118. spondylosis.mp.
119. spondylitis.mp.
120. spondylolisthesis.mp.
121. or/100-120
122. 88 and 121
123. exp neck/
124. exp cervical vertebrae/
125. Thoracic Vertebrae/
126. neck.mp.
127. (thoracic adj3 vertebrae).mp.
128. cervical.mp.
129. cervico*.mp.
130. 128 or 129
131. exp genital diseases, female/

132. genital disease*.mp.
133. exp *Uterus/
134. or/131-133
135. 130 not 134
136. (thoracic adj3 spine).mp.
137. cervical spine.mp.
138. 123 or 124 or 125 or 126 or 127 or 135 or 136 or 137
139. Intervertebral Disk/
140. (disc or discs).mp.
141. (disk or disks).mp.
142. 139 or 140 or 141
143. 138 and 142
144. herniat*.mp.
145. slipped.mp.
146. prolapse*.mp.
147. displace*.mp.
148. degenerat*.mp.
149. (bulge or bulged or bulging).mp.
150. 144 or 145 or 146 or 147 or 148 or 149
151. 143 and 150
152. intervertebral disk degeneration/ or intervertebral disk displacement/
153. intervertebral disk displacement.mp.
154. intervertebral disc displacement.mp.
155. intervertebral disk degeneration.mp.
156. intervertebral disc degeneration.mp.
157. 152 or 153 or 154 or 155 or 156
158. 138 and 157
159. 57 or 99 or 122 or 151 or 158
160. animals/ not (animals/ and humans/)
161. 159 not 160
162. exp *neoplasms/
163. exp *wounds, penetrating/
164. 162 or 163
165. 161 not 164
166. 29 and 165
167. guidelines as topic/
168. practice guidelines as topic/
169. guideline.pt.
170. practice guideline.pt.
171. (guideline? or guidance or recommendations).ti.
172. consensus.ti.
173. or/167-172
174. meta-analysis/
175. exp meta-analysis as topic/
176. (meta analy* or metaanaly* or met analy* or metanaly*).tw.
177. review literature as topic/

178. (collaborative research or collaborative review* or collaborative overview*).tw.
179. (integrative research or integrative review* or intergrative overview*).tw.
180. (quantitative adj3 (research or review* or overview*)).tw.
181. (research integration or research overview*).tw.
182. (systematic* adj3 (review* or overview*)).tw.
183. (methodologic* adj3 (review* or overview*)).tw.
184. exp technology assessment biomedical/
185. (hta or thas or technology assessment*).tw.
186. ((hand adj2 search*) or (manual* adj search*)).tw.
187. ((electronic adj database*) or (bibliographic* adj database*)).tw.
188. ((data adj2 abstract*) or (data adj2 extract*)).tw.
189. (analys* adj3 (pool or pooled or pooling)).tw.
190. mantel haenszel.tw.
191. (cohrane or pubmed or pub med or medline or embase or psycinfo or psyclit or psychinfo or psychlit or cinahl or science citation indes).ab.
192. or/174-191
193. 173 or 192
194. 166 and 193

# 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 |
| **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. | 3-5 |
| Objectives | 4 | Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS). | 4-5 |
| **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. | 5 |
| 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 |
| 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. | 6 |
| Search | 8 | Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated. | Appendix1 |
| 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). | 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. | 6-7 |
| Summary measures | 13 | State the principal summary measures (e.g., risk ratio, difference in means). | 8-9 |
| 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. | 8-9 |

# 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). | 8-9 |
| Additional analyses | 16 | Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified. | 8=9 |
| **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. | 9 |
| 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. | 9-10 |
| 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 |
| 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. | 10-12 |
| Synthesis of results | 21 | Present results of each meta-analysis done, including confidence intervals and measures of consistency. | 13 |
| Risk of bias across studies | 22 | Present results of any assessment of risk of bias across studies (see Item 15). | 10 |
| Additional analysis | 23 | Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]). | 13 |
| **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). | 14-15 |
| 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 |
| Conclusions | 26 | Provide a general interpretation of the results in the context of other evidence, and implications for future research. | 14-15 |
| **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. | 18 |

*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 the Global Rating of Change Scales in Patients with Neck Disorders: A Systematic Review with Meta-Analysis and Meta-Regression

SCHOLARONE™
Manuscripts

1 **Psychometric Properties of the Global Rating of Change Scales in Patients with Neck**

2 **Disorders: A Systematic Review with Meta-Analysis and Meta-Regression**

3 Pavlos Bobos[1], Joy C MacDermid[2], Goris Nazari[3], Rochelle Furtado[4] and CATWAD co-authors[5]

4

5 [1]Pavlos Bobos PT, PhD(c), (corresponding author) Doctoral Candidate, Western's Bone and Joint

6 Institute, Department of Health and Rehabilitation Sciences, Western University, Elborn College,

7 1201 Western Road, N6G 1H1, London, Ontario, Dalla Lana School of Public Health, Institute of

8 Health Policy Management and Evaluation, Department of Clinical Epidemiology and Health Care

9 Research, University of Toronto, Canada, (pbobos@uwo.ca), tel: +1 519 661 2111 x88912

10 [2]Joy C MacDermid BScPT, PhD, Professor, Physical Therapy and Surgery, Western University,

11 London, ON and Co-director Clinical Research Lab, Hand and Upper Limb Centre, St. Joseph's

12 Health Centre, London, Ontario; Professor Rehabilitation Science McMaster University,

13 Hamilton, ON, Canada (jmacderm@uwo.ca)

14 [3]Goris Nazari PT, PhD(c) Doctoral Candidate, Western's Bone and Joint Institute, School of

15 Physical Therapy, Department of Health and Rehabilitation Sciences, Western University,

16 London, Ontario, Canada, (gnazari@uwo.ca)

17 [4]Rochelle Furtado MSc Western's Bone and Joint Institute, School of Physical Therapy,

18 Department of Health and Rehabilitation Sciences, Western University, London, Ontario, Canada,

19 (rfurtad5@uwo.ca)

20 [5]CATWAD: Michele Sterling m.sterling@uq.edu.au, Anne Söderlund anne.soderlund@mdh.se,

21 Michele Curatolo, curatolo@uw.edu, James M Elliott j-elliott@northwestern.edu, David Walton

22 dwalton5@uwo.ca, Helge Kasch helgkasc@rm.dk, Linda Carroll linda.carroll@ualberta.ca,

23 Hans Westergren Hans.Westergren@skane.se, Gwendolen Jull g.jull@uq.edu.au, Eva-Maj

24 Malmström eva-maj.malmstrom@med.lu.se, Luke B Connelly l.connelly@uq.edu.au, Joy C

25 MacDermid jmacderm@uwo.ca, Mandy Nielsen mandy.nielsen@griffith.edu.au, Pierre Côté

26 pierre.cote@uoit.ca, Tonny Elmose Andersen tandersen@health.sdu.dk, Trudy Rebbeck

27 trudy.rebbeck@sydney.edu.au, Annick Maujean a.maujean@uq.edu.au, Sarah Robins

28 s.robins1@uq.edu.au, Kenneth Chen k.chen8@uq.edu.au, Julia Treleaven j.treleaven@uq.edu.au

**ABSTRACT**

31

32 **Objective:** The purpose of this systematic review was to critically appraise and synthesize the

33 psychometric properties of Global Rating of Change (GROC) scales for assessment of patients

34 with neck pain.

35 **Design:** Systematic review

36 **Data sources:** A search was performed in 4 databases (MEDLINE, EMBASE, CINAHL,

37 SCOPUS) until February 2019.

38 **Data extraction and synthesis:** Eligible articles were appraised using Consensus-based Standards

39 for the selection of health Measurement Instruments (COSMIN) checklist and the Quality

40 Appraisal for Clinical Measurement Research Reports Evaluation Form.

41 **Results:** The search obtained 16 eligible studies and included in total 1533 patients with neck pain.

42 Test-retest reliability of Global Perceived Effect (GPE) was very high (Intra-class correlation

43 coefficient (ICC) = 0.80 to 0.92) for patients with whiplash. Pooled data of Pearson's r indicated

44 that GROC scores were moderately correlated with neck disability change scores (0.53, 95% CI:

45 0.47 to 0.59). Pooled data of Spearman's correlations indicated that GROC scores were moderately

46 correlated with neck disability change scores (0.56, 95% CI: 0.41 to 0.68).

47 **Conclusions:** This study found excellent quality evidence of very good to excellent test-retest

48 reliability of GPE for patients with Whiplash Associated Disorders. Evidence from very good-to-

49 excellent quality studies found that GROC scores are moderately correlated to an external criterion

50 patient-reported outcome (PROM) measure evaluated pre-post treatment in patients with neck

51 pain. No studies were found that addressed the optimal form of GROC scales for patients with

52 neck disorders or compared the GROC to other options for single-item global assessment.

53 **Prospero registration number:** CRD 42018117874

54

**Strengths and limitations of this study**

- We rated the quality of individual studies and the overall risk of bias using two standardized approaches

- Our focus on neck pain increased the specificity of results but are not necessarily applicable to other musculoskeletal conditions

- Conceptual concerns about global ratings of change being affected by recall bias are not adequately addressed by psychometric evidence

- No studies addressing the optimal form of global rating were found.

**Introduction**

Neck pain is the 4th leading cause of disability and approximately half of adult the population with neck pain will experience a clinically important episode once in their lifetime. [1–3] The annual prevalence of neck pain it is estimated between 15% and 50%, with females having a higher prevalence rate than males. [2,3] Neck pain has been associated with many other comorbidities such as headaches, dizziness, anxiety, depression, back pain and arthralgias.[3–6] Several different methods for classifying neck pain have been described, using indicators such as duration (acute, sub-acute or chronic), degree of interference (low, moderate, severe) or most likely structure at fault (e.g. neuropathy vs. mechanical). [7]

As part of a patient-centric approach to care, clinicians will commonly evaluate response to intervention by asking the patient directly whether they feel better, worse, or the same since the prior encounter. While direct questioning can provide a qualitative indicator of change in status, many best practice guidelines endorse use of some form of quantified patient-reported outcome (PRO) as an adjunct to oral self-report. PROs are available to quantify several different constructs

3

79    in people with neck pain, including pain severity, disability and neck function. [8] Any PRO

80    intended to provide an estimate of change over time should be responsive to subtle shifts in the

81    patient's condition. To facilitate interpretation of change scores, a common property of many such

82    tools is the minimum clinically important difference (MCID), which is a change threshold that

83    corresponds to the minimum shift in scale values that most patients would indicate corresponds to

84    an important change in their overall condition. A well-recognized approach to establishing an

85    MCID for a PRO is to compare the magnitude of change against an anchor, most commonly a

86    Global Rating of Change (GROC) scale. These scales allow patients or study participants to

87    indicate whether their condition has gotten worse, better, or stayed the same and to quantify the

88    magnitude of that change. As they have been adopted as a sort of 'standard' against which change

89    in other tools is compared, the GROC can also be used on its own as an omnibus generic indicator

90    of change. [8]

91        Despite being accepted as a standard measure, there is considerable variation in how the

92    GROC has been constructed and implemented in research in neck pain. Some are 15 points, some

93    11 points, and others are 7 points. The common structure across these is the use of a middle '0'

94    score corresponding to 'no change', with negative values indicating magnitudes of worsening

95    while positive values indicate improvement.[9] Variations of the GROC (in name or structure)

96    include the "Global Perceived Effect", "Patient Global Impression of Change", "Transition

97    Ratings", and "Global Scale". [9]

98        A well-established component of health outcomes is having a tool with strong

99    psychometric properties of validity, reliability and responsiveness to be able to monitor change.

100    While recent research [8] has examined the psychometric properties of the most commonly

101    reported PROs for neck disorders, to date there has been no systematic review to summarize the

4

102 measurement properties of GROC scales themselves in patients with neck disorders. Therefore,

103 this systematic review aims to critically appraise and synthesize the psychometric properties of the

104 GROC scales in patients with neck disorders.

105

106 **METHODS**

107 *Patient and Public Involvement*

108 There was no patient or public involvement in the design or planning of this study.

109

110 *Study Design and Protocol Registration*

111 We conducted a systematic review to evaluate the psychometric properties of GROC scales in

112 patients with neck disorders. The protocol was registered in PROSPERO register database with

113 registration number: CRD 42018117874

114

115 *Eligibility Criteria*

116 We included studies in this systematic review if the following criteria were met [10–12]:

- 117 • Design: psychometric testing, randomized/ cohort studies

- 118 • Participants: > 50% of the study's patient population with neck conditions/disorders,

- 119 • Intervention/Comparison: studies that reported on the psychometric properties (reliability,
- 120 validity, responsiveness) of GROC, Global Perceived Effect (GPE) and Patient Global
- 121 Impression of Change (PGIC),

- 122 • Outcomes: GROC, GPE and PGIC

- 123 • Articles were written in English language only

5

124 Studies with no data on the GROC scales' psychometric properties, and conference

125 abstract/posters were excluded from this systematic review.

126

*Information Sources*

128 To identify studies on the psychometric properties (reliability, validity, responsiveness) of the

129 GROC, GPE and PGIC we searched the Medline, EMBASE, Scopus and CINAHL databases from

130 inception till February 2019, using a combination of keywords. Furthermore, we identified

131 additional studies by examining the reference list of each of the selected studies. The full list with

132 keyword strategy is presented in **APPENDIX 1**.

133

*Study Selection*

135 Two investigators (PB and GN) performed the systematic electronic searches independently in

136 each database. The same investigators then proceeded to identify and remove the duplicate studies.

137 In the next stage, we performed the independent screening of the titles and abstracts and any full-

138 text article marked as include or uncertain were obtained. In the final stage, the same two

139 independent authors performed the full text reviews independently to assess final article eligibility.

140 In case of disagreement, a third reviewer; the most experienced member (JM), facilitated a

141 consensus through discussion.

142

*Data Extraction*

144 The fourth author (RF) performed the data extractions. The extracted data were then cross-checked

145 by another author (PB). Data extraction included the author, year, study population/condition,

146 setting, sample size, age, properties evaluated, retest-interval, and the intervention protocol (if used

6

147  to assess responsiveness parameters). [13,14] For reliability estimates, Standard Error of

148  Measurement (SEM), Intra-class Correlation Coefficient (ICC), Minimal Detectable Change

149  (MDC) and 95% confidence intervals were extracted. [13,14] The ICC interpretation of ICC < 0.40

150  indicating poor, $0.40 \leq ICC < 0.75$ indicating fair-to-good and $ICC \geq 0.75$ indicating excellent

151  reliability were used as a common benchmark.[15] For validity estimates, correlation coefficient

152  (Pearson's/Spearman) and the 95% confidence intervals were extracted. [13,14] Evan's guidelines

153  to interpret the strength of the correlation was used which included: 0.00–0.19 "very weak", 0.20–

154  0.39 "weak", 0.40–0.59 "moderate", 0.60–0.79 "strong", and 0.80–1.00 "very strong". [16] For

155  responsiveness estimates, the Effect Size (ES), Standardized Response Mean (SRM), Clinically

156  Important Difference (CID), and/or Minimal Clinically Important Difference (MCID) including

157  the method of MCID estimation − Anchor-/Distribution-based methods, and 95% confidence

158  intervals were extracted. [13,14] To assist clinical decision making, standard benchmark scores of

159  trivial (< 0.20), small ($\geq 0.20$ to < 0.50), moderate ($\geq 0.50$ to < 0.80) or large ($\geq 0.80$), as proposed

160  by Cohen, were used. [17] When insufficient data were presented, PB contacted the authors by

161  email and requested further data.

162

163  *Consensus-based Standards for the selection of health Measurement Instruments (COSMIN)*

164  Consensus-based Standards for the selection of health Measurement Instruments (COSMIN)

165  assesses the risk of bias for the psychometric properties reported on a property-by-property basis.

166  A score for the risk of bias in estimates of psychometric properties was assessed by two authors

167  (PB) and (RF) using the new (COSMIN) checklist.[18] If disagreement was present a third person

168  (JM) assist in resolving the discrepancy. Each study was assessed by COSMIN on the 4-point scale

169  as "very good", "adequate", "doubtful" or "inadequate" for each of the checklist criteria for

7

170 relevant measurement properties (e.g. reliability, responsiveness, etc.). According to COSMIN,

171 when determining the overall score for each measurement property, the worst score counts method

172 was used wherein the lowest score for the checklist criteria of the relevant property was taken as

173 the overall score. [19] We then assessed the result of individual studies on a measurement property

174 against the updated criteria for good measurement properties. This involved the evaluation of

175 results of included studies as either sufficient (+), insufficient (−), or indeterminate (?). [18]

176

177 *Quality Appraisal for Clinical Measurement Research Reports Evaluation Form*

178 A summary score for the overall quality of individual studies was appraised independently by the

179 authors (PB) and (RF) using a structured clinical measurement specific appraisal tool. [13,14] In

180 case of disagreement a third person was consulted (JM) to resolve the conflict. The evaluation

181 criteria of this tool included twelve items: 1) Thorough literature review to define the research

182 question; 2) Specific inclusion/exclusion criteria; 3) Specific hypotheses; 4) Appropriate scope of

183 psychometric properties; 5) Sample size; 6) Follow-up; 7) The authors referenced specific

184 procedures for administration, scoring, and interpretation of procedures; 8) Measurement

185 techniques were standardized; 9) Data were presented for each hypothesis; 10) Appropriate

186 statistics-point estimates; 11) Appropriate statistical error estimates; and 12) Valid conclusions

187 and recommendations. [13,14] An article's total score – quality - was calculated by the sum of

188 scores for each item, divided by the numbers of items and multiplied by 100%. [13,14] Overall,

189 the quality summary of appraised articles range from (0%-30%) Poor, (31%-50%) Fair, (51%-

190 70%) Good, (71%-90%) Very Good, and (>90%) Excellent. [13,14]

191

192 *Synthesis of Results*

193    A qualitative synthesis was conducted to report findings on test-retest reliability statistics. A meta-

194    analysis of Pearson's and Spearman's correlation was performed in R (version 3.6.1) with

195    metaphor package.[20] The meta-analyses were conducted using a random effect model and the

196    correlation coefficients were converted to z values. Heterogeneity was deemed substantial if $I^2$

197    values were more than 50%. [21] A Meta-regression was planned to explore the sources of

198    unexplained heterogeneity by considering the following factors: a. neck pain with or without

199    radicular symptoms, b. acute or chronic, c. age and d. sex. Forest plots were created using means

200    and 95% confidence intervals for correlation coefficients. We summarize the main results of the

201    included articles based on the neck disorders, reported psychometric estimate and the study quality

202    ratings.

203

204    **RESULTS**

205    *Study Selection*

206    Our search yielded 8,837 articles. After removal of duplicates, 6,027 studies remained and were

207    screened using their title and abstract; leaving 29 articles selected for full-text review. Of these, 16

208    studies were considered eligible. [22,23,24–31,32–37] The flow of the study selection process is

209    presented in **Figure 1.**

210

211    *Study Characteristics*

212    The 16 eligible studies were conducted between 2006 and 2017 and included 1533 participants

213    with neck pain/disorders (mean of 96 participants per study). [22,23,24–31,32,34–37,] Study size

214    ranged from 29 to 200 participants. A summary description of all the studies included is displayed

215    in **Table 1.** Concurrent validity was evaluated in 14 studies by comparing the difference of pain

216 intensity, disability and function scores with the score of GROC scales. Two studies [26,31]

217 examined the test-retest reliability of a 7-point and an 11-point GPE scale for patients with

218 whiplash-associated disorders (WAD). One study [24] examined whether occurrences of within-

219 and between-session changes were significantly associated with functional outcomes, pain, and

220 self-report of recovery in patients at discharge who were treated with manual therapy for

221 mechanical neck pain.

222

223 *COSMIN Risk of Bias rating and Quality appraisal of the Included Studies*

224 Regarding the risk of bias, all studies were rated as very good (**Table 2**). The quality of the studies

225 ranged from 88% to 96% (**Table 3**). The most common flaws were 1) lack of/inadequate sample

226 size calculations, 2) missing data (i.e. inadequate follow up), and 3) inconsistencies between the

227 data presented and hypothesis stated.

228

229 *Reported GROC scales*

230 The most commonly reported GROC scale (n=6 studies) was a 15-point scale with the most

231 frequent anchors being "-7 (a very great deal worse) to zero (about the same) to +7 (a very great

232 deal better)". A 7-point scale was reported in 5 studies, 11- and 5-point scales were reported in 2

233 studies and a 9-point scale in one study. The anchors in those scales varied greatly and are

234 presented in Table 1. Only 6 studies [26,31–33,35,36] reported full detail regarding the specific

235 questions asked of the patients with neck disorder when a GROC scale was administered. Those

236 questions that were reported are presented in **Box 1.**

237

238

239 *Reliability Measures*

240 Two studies were included that examined test-retest reliability of GPE for patients with WAD.

241 Kamper et al. (2010) [26]  examined the [time interval] test-retest reliability of an 11-point GPE

242 scale in 134 patients with chronic WAD and reported an Intra-class Correlation Coefficient (ICC)

243 of 0.99 (95% CI 0.99 to 0.99) at baseline, 0.96 (0.95 to 0.97) at 6 weeks, and 0.92 (0.89 to 0.94)

244 at 12 months (**Table 4**). Ngo et al. (2010) assessed the test-retest reliability of a 7-point scale of

245 GPE in patients with acute WAD at 3 to 5 days. [31] The ICC and 95% confidence intervals (CI)

246 were used to determine the test–retest reliability of the two versions of the perceived recovery

247 questions using their original seven-item responses. Ngo et al. also computed weighted kappa

248 coefficients and 95% CI using quadratic weights to determine whether the distribution of responses

249 influenced the reliability as measured by the ICC. An ICC for general recovery of 0.70 (0.60 to

250 0.80) and an ICC for neck pain questions of 0.80 (0.72 to 0.87) were found. A weighted Kappa

251 was also calculated (Kappa = 0.70 (0.42 to 0.98)) at six weeks for general recovery and at six

252 weeks Kappa = 0.80 (0.51 to 1.0) for neck pain questions (**Table 4**).

253

254 *Validity Measures*

255 We found 14 studies that examined concurrent validity measures between GROC and another

256 PRO.[22,23,25,27–30,32,34,35,36–38] Correlations of Pearson's and Spearman's coefficients

257 between GROC and another PRO were ranging from very weak to very strong correlations. The

258 validity measures are presented and summarized in Table 5.

259

260

11

261 *Meta-Analysis and Meta-Regression of Correlations between Disability change scores and GROC*

262 *scores*

263 Five studies [23,25,34,37,38] of very good-to-excellent quality reported the Pearson correlation

264 coefficients between neck disability change scores and the GROC scores and were pooled together.

265 We found that GROC was positively correlated with disability change scores (r = 0.53, 95% CI:

266 0.47 to 0.59, $I^2$ = 0%). Six studies [27–30,32,36] of very good-to-excellent quality reported the

267 Spearman correlation coefficients between neck disability changes scores and the GROC scores

268 and were pooled together. We found that GROC was moderately correlated with disability change

269 scores (rho = 0.56, 95% CI: 0.41 to 0.68, $I^2$= 85%). The forest plots with correlation coefficients

270 with 95% CIs are presented in Figure 2-3. Our meta-regression showed that age was found as a

271 significant factor in influencing Fisher's Z scores (β = -0.034, 95% CI -0.05 to -0.01, p = 0.001).

272 The model explained 68% of the variance ($R^2$ = 0.68) (Figure 4).

273

274 *Area under the curve (AUC) – Sensitivity and Specificity*

275 Cook et al. [24] found that between-session NPRS- pain changes were associated with greater than

276 3-point change on the GROC at 96-hours (AUC=0.76). The pain change associated with GROC

277 was more specific (Specificity=79.2%, range: 62.2 - 91.1) than sensitive (Sensitivity=65.6%,

278 range: 57.9 to 74.6). Those with a 36.7% between-sessions change in pain were also 7.3 times

279 more likely to report an improvement of greater than 3 points change on the GROC than those

280 who did not achieve a 36.7% change in pain (**Table 4**).

281

282 **DISCUSSION**

283    This review has synthesized the current research from 16 studies that aimed to evaluate the

284    psychometric properties of GROC scales for patients with neck disorders, with the goal to provide

285    evidence for clinicians and researchers concerning its use within clinical practice and research.

286    From the 16 included studies, only 2 studies [26,31] reported test-retest reliability statistics of the

287    7- and 11-points item GPE scales for patients with WAD only. We were able to pool data from 12

288    studies regarding concurrent validity of GROC scales and neck disability change scores at one

289    time point after the interventions. Themes influencing interpretation of the GROC were explored

290    in a study [33] that evaluated the factors that contribute to how patients respond to a question on

291    global perceived effect. This study found that treatment process, biomechanical performance, self-

292    efficacy and the nature of the condition may influence the responses on global perceived effect,

293    which is consistent with what we would expect for patients with neck pain. This suggests that

294    change is a complex multifactorial global concept. A strength of GROC is that it is intended as a

295    global assessment, and it can be assumed that it reflects the aspects of change important to the

296    individual patient.

297    Reliability can be defined as the degree to which a measure produces consecutive results

298    with the least amount of random error when the status of the population remains unchanged. The

299    reliability of GPE displayed an excellent test-retest reliability of ICC>0.90 over an interval of 6

300    weeks and 12 months for patients with WAD. Conducting an assessment with a long test-retest

301    interval (e.g. 12 months), can provide challenges as there is higher risk of individuals with WAD

302    being symptomatically unstable.[9] Determining if patients are symptomatically-stable can be

303    achieved by administering another PRO such as the Single Assessment Numeric Evaluation

304    (SANE)[39], however, the 7- and 11- points GPE scales still demonstrated good stability properties

13

305 at long test intervals (i.e., of 6 weeks and 12 months).[26] Therefore, the measurements of the

306 reliability parameters of the GPE may be very useful during longer test intervals in clinical trials.

307 The psychometric property of validity is defined as the degree to which a PRO measures

308 what it is intended to measure. Pooled data from 11 studies overall suggest that post-treatment

309 changes of on validated disability outcome measures were moderately (Pearson's r = 0.51, 95%

310 CI: 0.43 to 0.58; Spearman's rho = 0.56, 95% CI: 0.41 to 0.68) correlated to change in perceived

311 effect) (Figure 2-3). This finding suggests that GROC scores taken at one point in time were related

312 to scores in pain and disability in patients with neck disorders, as measured by standardized

313 measures taken at 2 points in time. We identified one study [24] that found a 36.7% change in pain

314 for within- and between- session changes was associated with a 50% reduction in the NDI and an

315 improvement of >3 points on a 15-points GROC scale for patients with neck pain. This quantified

316 predictive change value may have clinical utility for use in clinical practice.

317 Previous studies [9,40] have indicated serious concerns about the conceptual validity of the

318 global rating of change. The review by Kamper et al.[9]  clearly showed that GROC was related

319 to final status more than change and was least related to baseline health status. This result

320 undermines the premise of what the global rating of change actually measures. For this reason, we

321 conclude that the 0.50 pooled correlation across 12 studies between the GROC and other PROM

322 change scores (e.g. NDI scores) may reflect a relationship between follow-up status and change

323 rather than supporting the contention that GROC actually measures change. This would also

324 explain why only 25% of the variation in GROC change scores was explained by changes scores

325 from a PROM change score measured at 2 points in time. In all studies, participants completed the

326 GROC scale at one time point after the intervention, and hence recall bias is a cause for concern.

327 However, another potential factor for moderate correlations is that the PROMs that have been used

14

328   as the comparator with GROC scores may not reflect priorities that are important to patients. That

329   is, the field has largely been driven by assumptions that the GROC is a 'gold standard' for

330   evaluating true change in a respondent's condition or status, and that all items on the comparator

331   PROM are of equal importance to all people with that condition. The work presented herein

332   challenges the valorization of the GROC as a gold standard for change, and prior work has

333   challenged the notions that all PROM items are equally important.[9,41,42] It is therefore possible

334   that the very constructs being evaluated require greater critical discourse before authors can say,

335   with confidence, that one scale functions well or poorly based on its associations with another

336   scale. Since no studies compared a retrospective global assessment of the GROC to pre-post single

337   item global PROM e.g. the SANE, we do not know the extent to which these two factors

338   contributed to moderate correlation.

339       A unique aspect of this study was that it focused on global rating of change scales in a neck

340   pain patient population. Our study appraisal suggests that future studies concerning GROC should

341   include adequate sample sizes, maintain a rigorous follow up and report appropriate statistical error

342   estimates, since these were often inadequate. Various critical appraisal tools exist, and the

343   perspectives and ratings may differ across instruments. We used 2 different critical appraisal tools

344   to evaluate quality from 2 perspectives. The COSMIN risk of bias assessments reflects the level

345   of confidence in the conclusions and pooled estimates. The quality appraisal tool focuses on design

346   issues in the studies and reflects gaps in research designs that should be considered in interpretation

347   of current research and improved in future studies. Substantial heterogeneity was detected

348   ($I^2 > 50\%$) in pooled Spearman's correlation coefficients which is a concern when pooling data. Our

349   univariate meta-regression analysis indicated that age across the studies explained 68% of the

350   variance (**Figure 4**). Other factors such as type of neck pain (with or without radicular symptoms),

15

351 acute or chronic and sex did not explain the remaining heterogeneity (not statically significant).

352 Furthermore, the scope of our literature search was focused on identifying full-text papers written

353 in English only.

354 While this study included 16 studies, only 2 of these reported reliability statistics for GROC

355 scales for patients with chronic WAD. Therefore, the applicability of our study is mostly limited

356 to patients with chronic WAD. For validity measurements, GROC scales were mostly investigated

357 by correlation analyses to evaluate the external responsiveness of another PRO measure over a

358 specific time point. From our meta-analysis, we can be confident that the GROC scores were

359 moderately correlated with neck disability change scores. However, more robust psychometric

360 design studies to test the measurement properties of GROC scales as the primary outcome of

361 investigation are highly needed. Future studies should aim to test to what extent the different range

362 of items (e.g. 7-point scale vs 11-point scale), the anchors (e.g. much worse vs much better) may

363 affect the measurement properties of GROC scales for patients with neck disorders. Also, it is

364 important to indicate that most outcome measures are ordinal and assume that additive scores of ordinal

365 items can be treated as interval level. This potentially could lead to scaling problems even in the face of

366 strong psychometric properties. The main protection we have is to create new scales or retrofit existing

367 scales based on Rasch analysis.

368

369 **CONCLUSIONS**

370 This study found excellent quality evidence of very good to excellent test-retest reliability of GPE

371 for patients with WAD. Evidence of very good to excellent quality studies found that GROC scores

372 are moderately correlated to an external criterion PROM measure measured pre-post treatment in

373 patients with neck disorders. Studies addressing the optimal form of GROC scales for patients with

16

374    neck disorders or comparing the GROC to other options for single-item global assessment of

375    change were not found.

376

**Authors' contributions**

378    PB contributed significantly to conception and design of the study, data extraction, critical

379    appraisal, interpretation of data and drafting of the manuscript. GN, and RF were involved in

380    literature search, critical appraisal and interpretation of data and drafting. GN was involved in

381    critical appraisal and drafting. JM was also involved in the conception and design of the study,

382    drafting, and revised the manuscript for important intellectual content. JM and CATWAD were

383    involved in the drafting and review of the manuscript. All authors have given their final approval

384    on the manuscript to be published

385

**Declarations**

387    **Ethics approval and consent to participate**

388    Not applicable

389    **Consent for publication**

390    Not applicable

391    **Availability of data and material**

392    Data sharing is not applicable to this article as no datasets were generated or analyzed during the

393    current study

17

397    **Competing Interest Statement**

398    None to report

399

400

401    **References**

402    1    Murray CJL, Abraham J, Ali MK, *et al.* The State of US health, 1990-2010: Burden of diseases,

403         injuries, and risk factors. *JAMA - J Am Med Assoc* Published Online First: 2013.

404         doi:10.1001/jama.2013.13805

405    2    Fejer R, Kyvik KO, Hartvigsen J. The prevalence of neck pain in the world population: A

406         systematic critical review of the literature. Eur. Spine J. 2006. doi:10.1007/s00586-004-0864-4

407    3    Hogg-Johnson S, van der Velde G, Carroll LJ, *et al.* The Burden and Determinants of Neck Pain in

408         the General Population. Results of the Bone and Joint Decade 2000-2010 Task Force on Neck

409         Pain and Its Associated Disorders. *J Manipulative Physiol Ther* Published Online First: 2009.

410         doi:10.1016/j.jmpt.2008.11.010

411    4    Bobos P, Nazari G, Palimeris S, *et al.* The contribution of health and psychological factors in

412         patients with chronic neck pain and disability: A cross-sectional study. *J Clin Diagnostic Res*

413         2018;**12**:YC04–7. doi:10.7860/JCDR/2018/31284.11203

414    5    MacDermid JC, Walton DM, Bobos P, *et al.* A Qualitative Description of Chronic Neck Pain has

415         Implications for Outcome Assessment and Classification. *Open Orthop J* 2017;**10**:746–56.

416         doi:10.2174/1874325001610010746

417    6    Treleaven J. Sensorimotor disturbances in neck disorders affecting postural stability, head and eye

418         movement control-Part 2: Case studies. *Man Ther* 2008;**13**:266–75.

419         doi:10.1016/j.math.2007.11.002

420    7    Cohen SP. Epidemiology, diagnosis, and treatment of neck pain. *Mayo Clin Proc* 2015;**90**:284–99.

421         doi:10.1016/j.mayocp.2014.09.008

18

422   8   Bobos P, Macdermid JC, Walton DM, *et al.* Patient-reported outcome measures used for neck

423        disorders: An overview of systematic reviews. J. Orthop. Sports Phys. Ther. 2018;**48**:775–88.

424        doi:10.2519/jospt.2018.8131

425   9   Kamper SJ, Maher CG, Mackay G. Global Rating of Change Scales: A Review of Strengths and

426        Weaknesses and Considerations for Design. *J Man Manip Ther* 2009;**17**:163–70.

427        doi:10.1002/mus.21062

428   10  Nazari G, Bobos P, MacDermid JC, *et al.* The Effectiveness of Instrument-Assisted Soft Tissue

429        Mobilization in Athletes, Participants Without Extremity or Spinal Conditions, and Individuals

430        with Upper Extremity, Lower Extremity, and Spinal Conditions: A Systematic Review. *Arch Phys*

431        *Med Rehabil* Published Online First: February 2019. doi:10.1016/j.apmr.2019.01.017

432   11  Bobos P, Nazari G, Szekeres M, *et al.* The effectiveness of joint-protection programs on pain,

433        hand function, and grip strength levels in patients with hand arthritis: A systematic review and

434        meta-analysis. *J Hand Ther* 2018;**32**:194–211. doi:10.1016/j.jht.2018.09.012

435   12  Nazari G, Bobos P, MacDermid JC, *et al.* Psychometric properties of the Zephyr bioharness

436        device: A systematic review. *BMC Sports Sci Med Rehabil* 2018;**10**. doi:10.1186/s13102-018-

437        0094-4

438   13  Law MC, MacDermid J. *Evidence-based rehabilitation : a guide to practice*. Thorofare, NJ: :

439        Slack Incorporated 2014.

440   14  Roy JS, Desmeules F, MacDermid JC. Psychometric properties of presenteeism scales for

441        musculoskeletal disorders: A systematic review. *J Rehabil Med* Published Online First: 2011.

442        doi:10.2340/16501977-0643

443   15  Sańchez J. Rosner, B.: Fundamentals of Biostatistics, third edition. PWS-Kent, Boston 1990, xv,

444        655 pp., ISBN 0-534-91973-1. *Biometrical J* 1993;**35**:150. doi:10.1002/bimj.4710350205

445   16  Wuensch KL, Evans JD. Straightforward Statistics for the Behavioral Sciences. *J Am Stat Assoc*

446        Published Online First: 2006. doi:10.2307/2291607

447   17  Cohen J. Statistical power analysis for the behavioral sciences. Stat. Power Anal. Behav. Sci.

19

448      1988. doi:10.1234/12345678

449   18   Mokkink LB, de Vet HCW, Prinsen CAC, *et al.* COSMIN Risk of Bias checklist for systematic

450        reviews of Patient-Reported Outcome Measures. *Qual Life Res* Published Online First: 2018.

451        doi:10.1007/s11136-017-1765-4

452   19   Terwee CB, Mokkink LB, Knol DL, *et al.* Rating the methodological quality in systematic reviews

453        of studies on measurement properties : a scoring system for the COSMIN checklist. 2012;:651–7.

454        doi:10.1007/s11136-011-9960-1

455   20   Viechtbauer W. Conducting meta-analisys in R with metafor package. *J Stat Softw* 2010;**36**:1–48.

456   21   Higgins JPT, Thompson SG, Deeks JJ, *et al.* Measuring inconsistency in meta-analyses. *BMJ*

457        Published Online First: 2003. doi:10.1136/bmj.327.7414.557

458   22   Cleland J, Fritz J, Whitman J, *et al.* The reliability and construct validity of the Neck Disability

459        Index and Patient Specific Functional Scale. *Spine (Phila Pa 1976)* 2006;**31**:598–602.

460   23   Cleland JA, Childs JD, Whitman JM. Psychometric Properties of the Neck Disability Index and

461        Numeric Pain Rating Scale in Patients With Mechanical Neck Pain. *Arch Phys Med Rehabil*

462        2008;**89**:69–74. doi:10.1016/j.apmr.2007.08.126

463   24   Cook C, Lawrence J, Michalak K, *et al.* Is there preliminary value to a within- and/or between-

464        session change for determining short-term outcomes of manual therapy on mechanical neck pain?

465        *J Man Manip Ther* 2014;**22**:173–80. doi:10.1179/2042618614y.0000000071

466   25   Guzy G, Vernon H, Polczyk R, *et al.* Psychometric validation of the authorized Polish version of

467        the Neck Disability Index. *Disabil Rehabil* 2013;**35**:2132–7. doi:10.3109/09638288.2013.771706

468   26   Kamper SJ, Ostelo RWJG, Knol DL, *et al.* Global Perceived Effect scales provided reliable

469        assessments of health transition in people with musculoskeletal disorders, but ratings are strongly

470        influenced by current status. *J Clin Epidemiol* 2010;**63**:760-766.e1.

471        doi:10.1016/j.jclinepi.2009.09.009

472   27   Shaheen AAM, Omar MTA, Vernon H. Cross-cultural adaptation, reliability, and validity of the

473        arabic version of neck disability index in patients with neck pain. *Spine (Phila Pa 1976)*

20

474    2013;**38**:609–15. doi:10.1097/BRS.0b013e31828b2d09

475    28    Takeshita K, Hosono N, Kawaguchi Y, *et al.* Validity, reliability and responsiveness of the

476          Japanese version of the Neck Disability Index. *J Orthop Sci* 2013;**18**:14–21. doi:10.1007/s00776-

477          012-0304-y

478    29    Trouli MN, Vernon HT, Kakavelakis KN, *et al.* Translation of the Neck Disability Index and

479          validation of the Greek version in a sample of neck pain patients. *BMC Musculoskelet Disord*

480          2008;**9**:1–8. doi:10.1186/1471-2474-9-106

481    30    Tuttle N, Laakso L, Barrett R. Change in impairments in the first two treatments predicts outcome

482          in impairments, but not in activity limitations, in subacute neck pain: An observational study. *Aust*

483          *J Physiother* 2006;**52**:281–5. doi:10.1016/S0004-9514(06)70008-3

484    31    Ngo Trung, Stupar Maja, Coˆteˊ Pierre, Boyle Eleanor, Shearer Heather. A study of the test –

485          retest reliability of the self-perceived general recovery and self-perceived change in neck pain

486          questions in patients with recent whiplash-associated disorders. 2010;:957–62.

487          doi:10.1007/s00586-010-1289-x

488    32    Björklund M, Wiitavaara B, Heiden M. Responsiveness and minimal important change for the

489          ProFitMap-neck questionnaire and the Neck Disability Index in women with neck–shoulder pain.

490          *Qual Life Res* 2017;**26**:161–70. doi:10.1007/s11136-016-1373-8

491    33    Evans R, Bronfort G, Maiers M, *et al.* '" I know it " s changed '': a mixed-methods study of the

492          meaning of Global Perceived Effect in chronic neck pain patients. 2014;:888–97.

493          doi:10.1007/s00586-013-3149-y

494    34    Jorritsma W, Dijkstra PU, De Vries GE, *et al.* Detecting relevant changes and responsiveness of

495          Neck Pain and Disability Scale and Neck Disability Index. *Eur Spine J* 2012;**21**:2550–7.

496          doi:10.1007/s00586-012-2407-8

497    35    Monticone M, Frigau L, Vernon H, *et al.* Responsiveness and minimal important change of the

498          NeckPix© in subjects with chronic neck pain undergoing rehabilitation. *Eur Spine J*

499          2018;**27**:1324–31. doi:10.1007/s00586-017-5343-9

21

500   36   Monticone M, Ambrosini E, Vernon H, *et al.* Responsiveness and minimal important changes for

501        the Neck Disability Index and the Neck Pain Disability Scale in Italian subjects with chronic neck

502        pain. *Eur Spine J* 2015;**24**:2821–7. doi:10.1007/s00586-015-3785-5

503   37   Young BA, Walker MJ, Strunce JB, *et al.* Responsiveness of the Neck Disability Index in patients

504        with mechanical neck disorders. *Spine J* 2009;**9**:802–8. doi:10.1016/j.spinee.2009.06.002

505   38   Farooq MN, Mohseni-Bandpei MA, Gilani SA, *et al.* Urdu version of the neck disability index: A

506        reliability and validity study. *BMC Musculoskelet Disord* 2017;**18**:1–11. doi:10.1186/s12891-017-

507        1469-5

508   39   Williams GN, Gangel TJ, Arciero RA, *et al.* Comparison of the single assessment numeric

509        evaluation method and two shoulder rating scales. Outcomes measures after shoulder surgery. *Am

510        J Sports Med* 1999;**27**:214–21. doi:10.1177/03635465990270021701

511   40   Schmitt J, Abbott JH. Global Ratings of Change Do Not Accurately Reflect Functional Change

512        Over Time in Clinical Practice. *J Orthop Sport Phys Ther* 2015;**45**:106–11.

513        doi:10.2519/jospt.2015.5247

514   41   Chiarotto A, Ostelo RW, Boers M, *et al.* A systematic review highlights the need to investigate the

515        content validity of patient-reported outcome measures for physical functioning in patients with

516        low back pain. *J Clin Epidemiol* 2018;**95**:73–93. doi:10.1016/j.jclinepi.2017.11.005

517   42   Ailliet L, Knol DL, Rubinstein SM, *et al.* Definition of the construct to be measured is a

518        prerequisite for the assessment of validity. the Neck Disability Index as an example. *J Clin

519        Epidemiol* 2013;**66**:775-782.e2. doi:10.1016/j.jclinepi.2013.02.005

520

521

522

523

524

525

22

526    **Figure 1.** Flow diagram of included studies

527    **Figure 2**. Meta-analysis of Pearson's correlation coefficients between neck disability change scores and
528    GROC scores in patients with neck disorders based on 5 very good to excellent quality studies.

529    **Figure 3**. Meta-analysis of Spearman's correlation coefficients between neck disability change scores
530    and GROC scores in patients with neck disorders based on 6 very good to excellent quality studies.

531    **Figure 4**. Random effects univariate meta-regression between age and the Fisher's Z estimates. Each circle
532    represents a study and the size of the circle indicates the influence of that study on the model. The
533    regression prediction is illustrated by the straight line and the curved lines represent the 95% confidence
534    intervals. Age explained 68% of the variance in the model ($R^2$=0.68)

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

23

**Table 1**. Study Characteristics

| Study | Population | Setting | Sample Size | Properties Evaluated | GROC evaluated | Interval |
|-------|-----------|---------|-------------|----------------------|----------------|----------|
| Bjorklund et al (2017) | Women with non-specific neck-shoulder pain | Not specified | 104 | Validity (correlation)<br><br>Between NDI and GRoC | GRoC 7-points<br><br>1. Very much worse; 2. Much worse; 3. Minimally worse; 4. No change; 5. Minimally improved; 6. Much improved; 7. Very much improved. | GRoC scale administered only after intervention at one time point (1 week) |
| Cleland et al (2006) | Patients with cervical radiculopathy | Hospital | 38 | Validity (correlation)<br><br>Between NDI and GRoC<br><br>Between PSFS and GRoC | GRoC 15-points<br><br>-7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at follow up. Within a week over the period of 7 weeks. |
| Cleland et al. (2008) | Patients with neck pain only | 5 Outpatient physical therapy clinics | 137 | Validity (correlation)<br><br>Between NDI and GRoC<br><br>Between NPRS and GRoC | GRoC 15-points<br><br>-7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at follow up. Within a week |
| Cook et al (2014) | Patients with any neck pain | Academic locations in Northeast Ohio | 56 | ROC curves and AUC to measure sensitivity and specificity. Binomial logistic regression analysis was also calculated to determine overall effect. | GRoC 15-points<br><br>-7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | Baseline and at follow up 48- and 96-hours post baseline |
| Farooq et al. (2017) | Patients with neck pain | Physical therapy clinics | 106 | Validity (correlation)<br><br>Between NDI-U and GRoC | GRoC 15-points<br><br>-7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at three weeks after intervention |
| Guzy et al. (2013) | Patients with neck pain | Outpatient rehabilitation clinic | 95 | Validity (correlation)<br><br>Between NDI-P and GRoC | GRoC 7-points<br><br>'complete recovery'' over ''no change'' to ''my complaints are worse than ever'' | GRoC scale was completed at 2 weeks and at 4 weeks |
| Jorritsma et al. (2012) | Patients with chronic non-specific neck pain | Tertiary university center for rehabilitation | 76 | Validity (correlation)<br><br>Between NDI and GRoC<br><br>Between NPAD and GRoC | GPE 7-points<br><br>3 (completely recovered) to zero (no change) to -3 (worse than ever) | After completion of the program varying from 3 to 5 months patients filled the GPE |
| Kamper et al. (2010) | Patients with any whiplash-associated disorder. | Physical therapy clinics | 134 | Test-retest reliability | GPE 11-points<br><br>-5 (vastly worse) to zero (unchanged) to +5 (completely recovered) | Baseline, 6 weeks and 12 months |
| Monticone et al. 2017 | Patients with chronic neck pain | Outpatient Rehabilitation Unit | 153 | Validity (correlation)<br><br>Between NeckPix and GPE | GPE 5-points<br><br>(helped a lot = 1, helped = 2), one no change level (helped only a little = 3), and two worsening | At the end of treatment (8 weeks) and one year before follow-up |

24

| | | | | | | |
|---|---|---|---|---|---|---|
| | | | | | levels (did not help = 4, made things worse = 5) | |
| Monticone et al. 2015 | Patients with chronic neck pain | Outpatient Rehabilitation Unit | 200 | Validity (correlation)<br><br>Between NDI and GPE<br><br>Between NPDS and GPE | GPE 5-points<br><br>(helped a lot = 1, helped = 2), one no change level (helped only a little = 3), and two worsening levels (did not help = 4, made things worse = 5) | At the end of treatment 8 weeks |
| Ngo et al. (2010) | Patients with WAD. Most participants (69.6%) had grade II WAD. | Interviewed by person or by telephone in Ontario | 46 | Test-retest reliability | GPE 7-points<br><br>1. General recovery question<br><br>Completely better Much improved Slightly improved No change<br><br>Slightly worse Much worse<br><br>Worse than ever<br><br>2. Change in neck pain question:<br><br>very much better, better, slightly better, no change, slightly worse, worse, or very much worse | 3-5 days |
| Shaheen et al. (2015) | Patients with neck pain lasting more than 3 months | 3 primary health centers | 70 | Validity (correlation)<br><br>Between NDI-Ar and GRoC | GRoC 15-points<br><br>-7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | 1 week |
| Takeshita et al. (2014) | Patients with neck pain, cervical radiculopathy and/or cervical myelopathy | Variety of clinics and hospital settings | 130 | Validity (correlation)<br><br>Between NDI-J and GRoC | PGIC 7-points<br><br>much better, better, slightly better, unchanged, slightly worse, worse and much worse | Over 8 weeks |
| Trouli et al. (2008) | Patients with neck pain | Primary healthcare clinic | 68 | Validity (correlation)<br><br>Between NDI-Gr and GRoC | GRoC 15-points<br><br>-7 (a very great deal worse) to -1 (almost the same, hardly any worse at all) and from 7 (a very great deal better) to 1 (almost the same, hardly any better at all) | Within 2 months, but 1 week for test-retest |
| Tuttle et al. (2006) | Patients with neck pain for more than 2 weeks | Private physiotherapy clinics | 29 | Validity (correlation)<br><br>Between NDI and GPE<br><br>Between PSFS and GPE<br><br>Between VAS and GPE<br><br>Between ROM and GPE | GPE 11-points<br><br>−5 is vastly worse and +5 is completely recovered | 6 weeks |
| Young et al. (2009) | Patients presenting with mechanical neck pain | Outpatient physical therapy | 91 | Validity (correlation) | GRoC 15-points<br><br>-7 (''a very great deal worse'') to 0 (''about the same'') to +7 (''a | 3 weeks |

clinics.

very great deal better'')

556

26

558 **TABLE 2.** Summary of Psychometric Properties Reported in Studies and COSMIN Risk of Bias (RoB)
559 and Quality studies

| Study | Psychometric Properties Reported | COSMIN RoB | COSMIN Rating*§ (Criteria) | Quality of Studies** (QACMRR) |
|---|---|---|---|---|
| Bjorklund et al (2017) | Validity (correlation) | Very Good | ? | Excellent |
| Cleland et al (2006) | Validity (correlation) | Very Good | + | Excellent |
| Cleland et al. (2008) | Validity (correlation) | Very Good | - | Excellent |
| Cook et al (2014) | Sensitivity Specificity | Very Good Very Good | + | Excellent |
| Farooq et al. (2017) | Validity (correlation) | Very Good | + | Excellent |
| Guzy et al. (2013) | Validity (correlation) | Very Good | ? | Very good |
| Jorritsma et al. (2012) | Validity (correlation) | Very Good | ? | Excellent |
| Kamper et al. (2010) | Test-retest reliability | Very Good | + | Excellent |
| Monticone et al. (2017) | Validity (correlation) | Very Good | ? | Excellent |
| Monticone et al. (2015) | Validity (correlation | Very Good | ? | Excellent |
| Ngo et al. (2010) | Test-retest reliability | Very Good | + | Excellent |
| Shaheen et al. (2015) | Validity (correlation) | Very Good | ? | Excellent |
| Takeshita et al. (2014) | Validity (correlation) | Very Good | ? | Very good |
| Trouli et al. (2008) | Validity (correlation) | Very Good | + | Excellent |
| Tuttle et al. (2006) | Validity (correlation) | Very Good | ? | Excellent |
| Young et al. (2009) | Validity (correlation) | Very Good | ? | Excellent |

560 COSMIN, Consensus-based Standards for the Selection of health Measurement Instruments, Criteria for good measurement
561 properties: '+' sufficient; '-'insufficient; '?' indeterminate. §§ The grading for the quality of the evidence based on the modified
562 GRADE approach is not applicable. **Quality Appraisal for Clinical Measurement Research Reports Evaluation Form
563 (QACMRR).

564

565

566

567

568

569

570

571

572

27

573 **TABLE 3**. Quality Appraisal for Clinical Measurement Research Reports Evaluation Form

| Study | \multicolumn{12}{c}{Item Evaluation Criteria*} | Total (%) | Quality Summary |
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bjorklund et al (2017) | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Cleland et al. (2008) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Trouli et al. (2008) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Tuttle et al. (2006) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Kamper et al. (2010) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Cook et al (2014) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 92 | Excellent |
| Jorritsma et al. (2012) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Cleland et al (2006) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Monticone et al. (2017) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Monticone et al. (2015) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Ngo et al. (2010) | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 1 | 2 | 92 | Excellent |
| Shaheen et al. (2013) | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 92 | Excellent |
| Farooq et al. (2017) | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 92 | Excellent |
| Young et al. (2009) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Guzy et al. (2013) | 2 | 2 | 1 | 2 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 88 | Very good |
| Takeshita et al. (2014) | 2 | 2 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 88 | Very good |

574 *Item Evaluation Criteria: 1. Thorough literature review to define the research question; 2. Specific inclusion/exclusion

575 criteria; 3. Specific hypotheses; 4. Appropriate scope of psychometric properties; 5. Sample size; 6. Follow-up; 7. The

576 authors referenced specific procedures for administration, scoring, and interpretation of procedures; 8. Measurement

577 techniques were standardized; 9. Data were presented for each hypothesis; 10. Appropriate statistics-point estimates; 11.

578 Appropriate statistical error estimates; 12. Valid conclusions and clinical recommendations.*

28

579   *Total score = (sum of subtotals ÷ 24 × 100). If for a specific paper an item is deemed NA (Not Applicable), then, Total score*

580   *= (sum of subtotals ÷ (2 × number of Applicable items) × 100).*

581   *NA – Not Applicable. The subsections no. 6, asks for percentage of retention/follow up. This subsection only applies to*

582   *reliability test-retest studies*

583   *Quality Summary: Poor (0%-30%), Fair (31%-50%), Good (51%-70%), Very good (71%-90%), Excellent (>90%):*

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

29

608 **TABLE 4**. Summary of reliability properties of GRoC scales

| Study | Type of Reliability | Reliability Estimates | COSMIN | Quality of Studies |
|---|---|---|---|---|
| Kamper et al. (2010) | Test-retest | Intra-class correlation coefficients (ICC)<br>0.99 (0.99 – 0.99) – baseline<br>0.96 (0.95 – 0.97) – at six weeks<br>0.92 (0.89 – 0.94) at twelve months. | Very Good | Excellent |
| Ngo et al. (2010) | Test-retest | Intra-class correlation coefficients (ICC)<br>0.70 (0.60–0.80) – at six weeks (General recovery)<br>0.80 (0.72–0.87) – at six weeks (neck pain questions)<br><br>Weighted Kappa<br>0.70 (0.42–0.98) – at six weeks (General recovery)<br>0.80 (0.51–1.0) – at six weeks (neck pain questions)<br><br>Dichotomized response options for recovery (K statistics)<br>0.85 (0.64–1) when ''recovered'' was defined ''completely better'<br>0.81 (0.64–0.99) when defined as ''completely better'' or ''much improved<br><br>Dichotomized response options for change in neck pain questions (K statistics)<br>0.46 (0.20–0.74) when ''recovered'' was defined as ''very much better''<br>0.80 (0.62–0.99) when defined as ''very much better'' or ''better'<br><br>Recall questions (K statistics)<br>the kappa coefficient was 1 for participants who remembered their previous answers to the general recovery question; 0.88 (0.64–1) for those who did not remember and 0.50 (0.02– 0.98) for participants who were not asked the question.<br><br>The kappa coefficient was 1 for participants who remembered their previous answers to the change in neck pain question; 0.74 (0.41–1) for those who did not remember and 0.66 (0.22–1) for participants who were not asked the question. | Very Good | Excellent |

609

610

611

612

613

614

615

616

617

30

618 **TABLE 5**. Summary of validity properties of GRoC scales

| Study | Type of Reliability | Validity Estimates | COSMIN | Quality of Studies |
|---|---|---|---|---|
| Bjorklund et al (2017) | Spearman's correlation between the change scores of GRoC and ProFitMap-neck<br><br>GRoC and NDI | rho = 0.47, (p<0.05)<br>rho = 0.59, (p<0.05) | Very Good | Excellent |
| Cleland et al. (2006) | Correlations (Pearson r) between change scores NDI and GRoC<br>PSFS and GRoC | r = 0.19<br><br>r = 0.82 | Very Good | Excellent |
| Cleland et al. (2008) | Correlations (Pearson r) between change scores NDI and GRoC<br>NRS and GRoC | r = 0.58<br>r = 0.57 | Very Good | Excellent |
| Cook et al. (2014) | Receiver operator characteristics (ROC) Within-session change Between-session change<br><br>Between session change of Pain and GROC Sensitivity Specificity | AUC = 0.61<br>AUC = 0.76, >36.7% change in pain<br><br>Odds ratio = 7.3 (2.1, 24.7)<br>65.6% (57.9, 74.6)<br>79.2% (62.2, 91.1) | Very Good | Excellent |
| Farooq et al. (2017) | Correlations (Pearson r) NDI-U | r =0.50 | Very Good | Excellent |
| Guzy et al. (2013) | Correlations (Pearson r) NDI vs GROC | Two- week interval (r = -0.73)<br>Four-week interval (r = -0.56) | Very Good | Very good |
| Jorritsma et al. (2012) | Correlation between change scores of NPAD and GPE | r = 0.49 (95 % CI 0.30–0.64) | Very Good | Excellent |
| Monticone et al. (2017) | Correlations (Spearman) between change scores of the NeckPix© and GPE | rho = 0.69–0.82 | Very Good | Excellent |
| Monticone et al. (2015) | Correlation (Spearman) between change scores NDI-I and GPE<br>NDPS and GPE | rho = 0.71, p<0.01<br>rho = 0.59, p<0.01 | Very Good | Excellent |
| Shaheen et al. (2013) | Correlations (Spearman's) NDI-Ar and GROC | rho = 0.81, p<o.oo1 | Very Good | Excellent |
| Takeshita et al. (2014) | Correlations NDI and PGIC<br>NDI-J and PGIC | Spearman (rho)<br>rho = 0.47, p<o.oo1<br>rho = 0.59, p<o.oo1 | Very Good | Very good |
| Trouli et al. (2008) | Correlation (Spearman's) GROC vs Gr-NDI | rho = 0.30, p=0.02 | Very Good | Excellent |
| Tuttle et al. (2006) | Correlations (Spearman's) NDI vs GPE (post 1, minus pre-1)<br>NDI vs GPE (post 2, minus pre-1)<br>NDI vs GPE (post 2, minus pre-2)<br><br>PSFS vs GPE (post 1, minus pre-1)<br>PSFS vs GPE (post 2, minus pre-1)<br>PSFS vs GPE (post 2, minus pre-2) | rho = 0.17<br>rho = 0.01<br>rho = 0.03<br><br>rho = 0.06<br>rho = 0.03<br>rho = 0.03 | Very Good | Excellent |

31

| | | | | |
|---|---|---|---|---|
| | Pain Intensity (post 1, minus pre-1) | rho = 0.00 | | |
| | Pain Intensity (post 2, minus pre-1) | rho = 0.05 | | |
| | Pain Intensity (post 2, minus pre-2) | rho = 0.01 | | |
| | Total ROM (post 1, minus pre-1) | rho = 0.03 | | |
| | Total ROM (post 2, minus pre-1) | rho = 0.01 | | |
| | Total ROM (post 2, minus pre-2) | rho = 0.00 | | |
| Young et al. (2009) | Correlations (Pearson's) between change scores NDI and GRoC | r =0.52 (p<0.01) | Very Good | Excellent |
| Monticone et al. (2015) | Correlation (Spearman) between change scores NDI-I and GPE NDPS and GPE | rho = 0.71, p<0.01 rho = 0.59, p<0.01 | Very Good | Excellent |

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

32

641

642 **Box 1.** Questions of Global Rating of Change (GROC) scales

| Author | GROC item- scale | Patients with neck disorders were asked: |
|---|---|---|
| Bjorklund et al. (2017) | GROC 7-points | *"Compared to before the treatment of the study started, my overall status is now"* <br><br> *"Compared to before the treatment of the study started, my status regarding my neck–shoulder problem is now"* |
| Evans et al (2014) | GPE 9-points | *"Overall, how much has your neck pain changed since you started treatment in the study?"* |
| Kamper et al. (2010) | GPE 11-points | *"With respect to your whiplash injury how would you describe yourself now compared to immediately after your accident"* |
| Monticone et al. (2017) | GPE 5-points | *"Overall, how much did the treatment you received help your fear of movement due to current neck pain?* <br><br> *"Overall, how much did the treatment you delivered help your subject's fear of movement due to her/ his current neck pain?"* |
| Monticone et al. (2015) | GPE 5-points | *"Overall, how much did the treatment you received help your neck problem?"* |
| Ngo et al. (2010) | GPE 7-points | *"How well do you feel you are recovering from your injuries?"* <br><br> *"How do you feel your neck pain has changed since the injury?"* |

643

644

645

646

647

33

Figure 1. Flow diagram of included studies
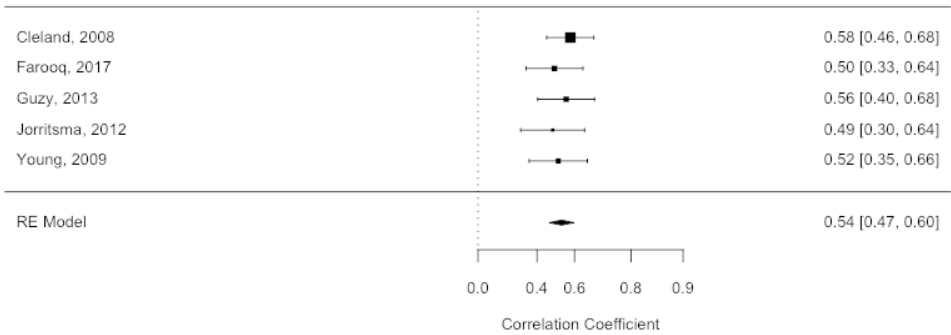
60x34mm (300 x 300 DPI)

Figure 2. Meta-analysis of Pearson's correlation coefficients between neck disability change scores and GROC scores in patients with neck disorders based on 5 very good to excellent quality studies.
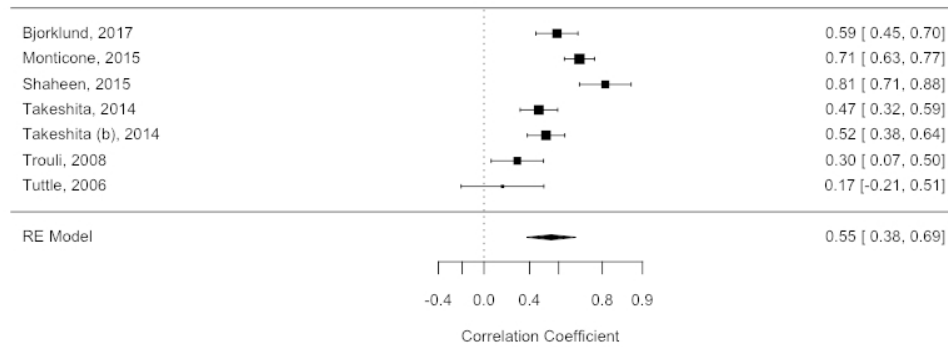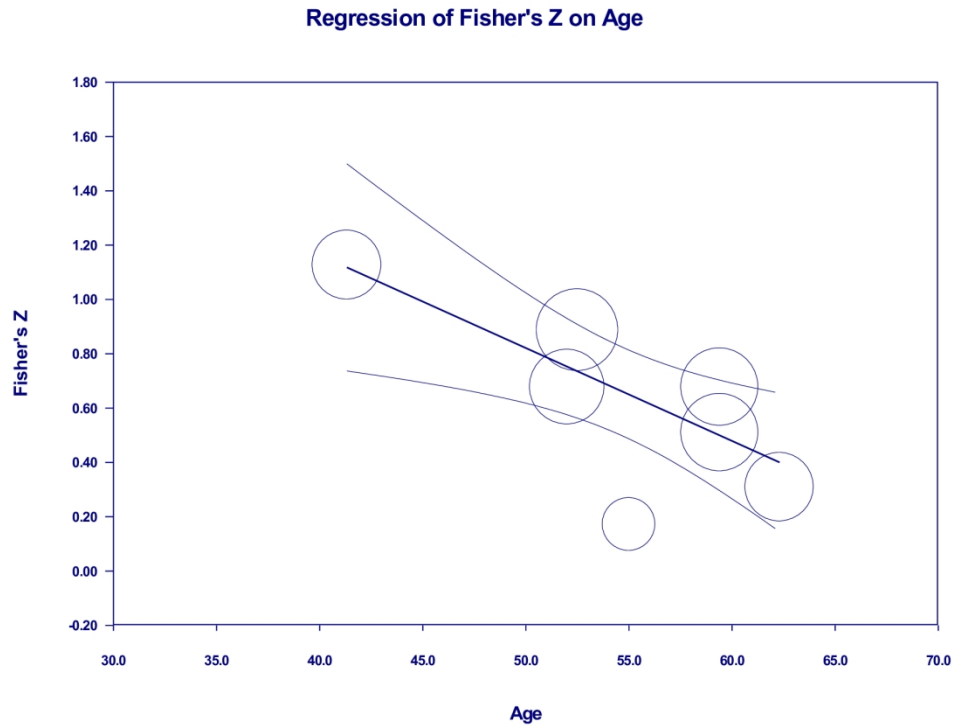
67x34mm (300 x 300 DPI)

Figure 3. Meta-analysis of Spearman's correlation coefficients between neck disability change scores and GROC scores in patients with neck disorders based on 6 very good to excellent quality studies.

67x34mm (300 x 300 DPI)

Figure 4. Random effects univariate meta-regression between age and the Fisher's Z estimates. Each circle represents a study and the size of the circle indicates the influence of that study on the model. The regression prediction is illustrated by the straight line and the curved lines represent the 95% confidence intervals. Age explained 68% of the variance in the model (R2=0.68)

160x118mm (300 x 300 DPI)

**Appendix 1**

**Search terms**

MEDLINE-OVID

1. exp "outcome and process assessment (health care)"/ or "outcome assessment (health care)"/ or treatment outcome/
2. outcome?.ti.
3. exp "Range of Motion, Articular"/
4. Pain Measurement/
5. exp disability evaluation/
6. "Recovery of Function"/
7. Questionnaires/
8. self-report.tw.
9. ((impairment or disability or function) adj2 (measure? or scale? or evaluation?)).tw.
10. range of motion.tw.
11. (strength adj2 (measure? or scale? or evaluation?)).tw.
12. (outcome? adj2 (measure* or scale? or indicator?)).tw.
13. or/1-12
14. "reproducibility of results"/
15. exp "Sensitivity and Specificity"/
16. reliability.mp.
17. validity.mp.
18. responsiveness.mp.
19. Psychometrics/
20. rasch.mp.
21. factor analysis, statistical/
22. factor analysis.tw.
23. differential functioning.mp.
24. (validity or validation).mp. [mp=title, original title, abstract, name of substance word, subject heading word, unique identifier]
25. (validity or validation).mp.
26. item difficulty.mp.
27. translation.tw.
28. or/14-27
29. 13 and 28
30. Neck Pain/
31. exp Brachial Plexus Neuropathies/
32. exp neck injuries/ or exp whiplash injuries/
33. cervical pain.mp.
34. neckache.mp.
35. whiplash.mp.
36. cervicodynia.mp.
37. cervicalgia.mp.
38. brachialgia.mp.
39. brachial neuritis.mp.

40. brachial neuralgia.mp.
41. neck pain.mp.
42. neck injur*.mp.
43. brachial plexus neuropath*.mp.
44. brachial plexus neuritis.mp.
45. thoracic outlet syndrome/ or cervical rib syndrome/
46. Torticollis/
47. exp brachial plexus neuropathies/ or exp brachial plexus neuritis/
48. cervico brachial neuralgia.ti,ab.
49. cervicobrachial neuralgia.ti,ab.
50. (monoradicul* or monoradicl*).tw.
51. or/30-50
52. exp headache/ and cervic*.tw.
53. exp genital diseases, female/
54. genital disease*.mp.
55. or/53-54
56. 52 not 55
57. 51 or 56
58. neck/
59. neck muscles/
60. exp cervical plexus/
61. exp cervical vertebrae/
62. atlanto-axial joint/
63. atlanto-occipital joint/
64. Cervical Atlas/
65. spinal nerve roots/
66. exp brachial plexus/
67. (odontoid* or cervical or occip* or atlant*).tw.
68. axis/ or odontoid process/
69. Thoracic Vertebrae/
70. cervical vertebrae.mp.
71. cervical plexus.mp.
72. cervical spine.mp.
73. (neck adj3 muscles).mp.
74. (brachial adj3 plexus).mp.
75. (thoracic adj3 vertebrae).mp.
76. neck.mp.
77. (thoracic adj3 spine).mp.
78. (thoracic adj3 outlet).mp.
79. trapezius.mp.
80. cervical.mp.
81. cervico*.mp.
82. 80 or 81
83. exp genital diseases, female/
84. genital disease*.mp.
85. exp *Uterus/

86. 83 or 84 or 85
87. 82 not 86
88. 58 or 59 or 60 or 61 or 62 or 63 or 64 or 65 or 66 or 67 or 68 or 69 or 70 or 71 or 72 or 73 or 74 or 75 or 76 or 77 or 78 or 79 or 87
89. exp pain/
90. exp injuries/
91. pain.mp.
92. ache.mp.
93. sore.mp.
94. stiff.mp.
95. discomfort.mp.
96. injur*.mp.
97. neuropath*.mp.
98. or/89-97
99. 88 and 98
100. Radiculopathy/
101. exp temporomandibular joint disorders/ or exp temporomandibular joint dysfunction syndrome/
102. myofascial pain syndromes/
103. exp "Sprains and Strains"/
104. exp Spinal Osteophytosis/
105. exp Neuritis/
106. Polyradiculopathy/
107. exp Arthritis/
108. Fibromyalgia/
109. spondylitis/ or discitis/
110. spondylosis/ or spondylolysis/ or spondylolisthesis/
111. radiculopathy.mp.
112. radiculitis.mp.
113. temporomandibular.mp.
114. myofascial pain syndrome*.mp.
115. thoracic outlet syndrome*.mp.
116. spinal osteophytosis.mp.
117. neuritis.mp.
118. spondylosis.mp.
119. spondylitis.mp.
120. spondylolisthesis.mp.
121. or/100-120
122. 88 and 121
123. exp neck/
124. exp cervical vertebrae/
125. Thoracic Vertebrae/
126. neck.mp.
127. (thoracic adj3 vertebrae).mp.
128. cervical.mp.
129. cervico*.mp.

130. 128 or 129
131. exp genital diseases, female/
132. genital disease*.mp.
133. exp *Uterus/
134. or/131-133
135. 130 not 134
136. (thoracic adj3 spine).mp.
137. cervical spine.mp.
138. 123 or 124 or 125 or 126 or 127 or 135 or 136 or 137
139. Intervertebral Disk/
140. (disc or discs).mp.
141. (disk or disks).mp.
142. 139 or 140 or 141
143. 138 and 142
144. herniat*.mp.
145. slipped.mp.
146. prolapse*.mp.
147. displace*.mp.
148. degenerat*.mp.
149. (bulge or bulged or bulging).mp.
150. 144 or 145 or 146 or 147 or 148 or 149
151. 143 and 150
152. intervertebral disk degeneration/ or intervertebral disk displacement/
153. intervertebral disk displacement.mp.
154. intervertebral disc displacement.mp.
155. intervertebral disk degeneration.mp.
156. intervertebral disc degeneration.mp.
157. 152 or 153 or 154 or 155 or 156
158. 138 and 157
159. 57 or 99 or 122 or 151 or 158
160. animals/ not (animals/ and humans/)
161. 159 not 160
162. exp *neoplasms/
163. exp *wounds, penetrating/
164. 162 or 163
165. 161 not 164
166. 29 and 165
167. guidelines as topic/
168. practice guidelines as topic/
169. guideline.pt.
170. practice guideline.pt.
171. (guideline? or guidance or recommendations).ti.
172. consensus.ti.
173. or/167-172
174. meta-analysis/
175. exp meta-analysis as topic/

176. (meta analy* or metaanaly* or met analy* or metanaly*).tw.

177. review literature as topic/

178. (collaborative research or collaborative review* or collaborative overview*).tw.

179. (integrative research or integrative review* or intergrative overview*).tw.

180. (quantitative adj3 (research or review* or overview*)).tw.

181. (research integration or research overview*).tw.

182. (systematic* adj3 (review* or overview*)).tw.

183. (methodologic* adj3 (review* or overview*)).tw.

184. exp technology assessment biomedical/

185. (hta or thas or technology assessment*).tw.

186. ((hand adj2 search*) or (manual* adj search*)).tw.

187. ((electronic adj database*) or (bibliographic* adj database*)).tw.

188. ((data adj2 abstract*) or (data adj2 extract*)).tw.

189. (analys* adj3 (pool or pooled or pooling)).tw.

190. mantel haenszel.tw.

191. (cohrane or pubmed or pub med or medline or embase or psycinfo or psyclit or psychinfo or psychlit or cinahl or science citation indes).ab.

192. or/174-191

193. 173 or 192

194. 166 and 193

**Quality Appraisal for Clinical Measurement Research Reports**

**Evaluation Form**

Authors: _____ Year: _____ Rater: _____

*Use this form to rate the quality of a clinical measurement study. To decide which score to provide for each item on your quality checklist, pick the descriptor that sounds <u>most</u> like what was reported in the study you are evaluating. Items rank descriptors are provided in the guide. (Forms and guides to extract study data for evidence synthesis are available from developer at macderj@mcmaster.ca)*

| Evaluation criteria | Score | | |
|---|---|---|---|
| **Study question** | 2 | 1 | 0 |
| 1. Was the relevant background work cited to define what is currently known about the measurement properties of measures under study, and the potential contributions of the current research question to informing that knowledge base? | | | |
| **Study Design** | | | |
| 2. Were appropriate inclusion/exclusion criteria defined? | | | |
| 3. Were specific clinical measurement questions/hypotheses identified? | | | |
| 4. Was an appropriate scope of measurement properties considered? | | | |
| 5. Was an appropriate sample size used? | | | |
| 6. Was appropriate retention/follow-up obtained? (for studies involving retesting; otherwise  n/a) | | | |
| **Measurements** | | | |
| 7. Were specific descriptions provided of the measure under study and the method(s) used to administer it? | | | |
| 8.  Were standardized procedures used to administer all study measures in a manner that minimized potential sources of error/bias (including the study measure and its comparators)? | | | |
| **Analyses** | | | |

| | | | |
|---|---|---|---|
| 9. Were analyses conducted for each specific hypothesis or purpose? | | | |
| 10. Were appropriate statistical tests performed to obtain point estimates of the measurement properties? | | | |
| 11. Were appropriate ancillary analyses done to quantify the confidence in the estimates of the clinical measurement property (Precision/Confidence intervals; benchmark comparisons/ROC curves, alternate forms of analysis like SEM/MID, etc.)? | | | |
| **Recommendations** | | | |
| 12. Were clear, specific and accurate conclusions made about the clinical measurement properties; that were associated with appropriate clinical measurement recommendations and supported by the study objectives, analysis and results? | | | |
| **Subtotals** (of columns 1 and 2) | | | |
| **Total score** (sum of subtotals/24*100); if for a specific paper or topic an item is deemed inappropriate then you can sum of items/2*number of items *100 | | | |

© MacDermid 2011

## Quality Appraisal of a Clinical Measurement Study

### Interpretation Guide

To decide which score to provide for each item on your quality checklist, read the following descriptors. Pick the descriptor that sounds _most_ like the study you were evaluating with respect to a given item. If there is no documentation about any specific aspect of an item; then you must evaluate assuming that it was not done. Given the diversity in clinical measurement properties and design options, the evaluator has to make judgments using the criteria below and extend the principles to specific aspects that may not be covered in these brief exemplars. In many cases, the study will not look exactly like the descriptor so there will be some interpretation as to which level of optimal methods for clinical measurement studies have been achieved. In such cases, the evaluator can use the general approach that if this study research design and conduct is consistent with best practice (score=2); is acceptable but suboptimal (score=1); is not done/documented, substantially inadequate or inappropriate (score=0).

| | **Descriptors** | |
|---|---|---|
| **Study question** | | |
| Score | | |
| 1 | 2 | The authors:<br><br>- performed a thorough literature review indicating what is currently known, and not known, about the clinical measurement properties of the instruments or tests under study<br>- presented a critical, and unbiased view of what is known about the current measurement properties<br>- indicated how the current research question fills a gap in the current knowledge base<br>- established a research question based on the above. |
| | 1 | All of the above criteria were not fulfilled, but a sound rationale was provided for the research question. |
| | 0 | A foundation for the current research question was not clear; and the rationale was not founded on previous literature. |
| **Study design** | | |

| 2 | 2 | Specific inclusion/exclusion criteria for the study were defined, that described the patients enrolled. The subjects were described in terms of health condition/demographics, key relevant outcome mediators and the recruitment context (setting). |
|---|---|---|
| | 1 | Some information on participants and place is provided (not all of above). For example, age/sex/diagnosis and the name or type of the practice is listed; but no additional information. |
| | 0 | No information on type of clinical settings or study participants is provided (other than number/mean age). |
| 3 | 2 | Specific hypotheses or research questions are provided. The stated study purpose provides specific research questions or hypotheses that indicate which specific measurement properties will be evaluated. This should include the specific type of reliability (intra/inter-rater or test-retest) being tested or the type of validity (construct/criterion/content; longitudinal/concurrent; convergent/divergent) being tested. A prior hypothesis should describe the level of reliability expected; and for validity, expected relationships (strength of associations) or constructs. |
| | 1 | The types of reliability and validity being tested were apparent in the methods/title, but clear and specific research questions or hypotheses were not specified. |
| | 0 | Specific types of reliability or validity under evaluation were not clearly defined nor were specific hypotheses on reliability and validity stated. ("*The purpose of this study was to investigate the reliability and validity of...*" can be rated as zero if no further detail on the types of reliability and validity or the nature of specific hypotheses is stated). |
| 4 | 2 | An appropriate scope of clinical measurement properties would be indicated by<br><br>1. A detailed focus on reliability that included multiple forms of reliability (at least two of – intra-rater, inter-rater, test retest); as well as both relative and absolute reliability (e.g., ICCs and SEM/MID or limits of agreement)<br>2. A detailed focus on validity that included multiple forms of validity (content (judgmental); structured (e.g., expert review/survey, qualitative interviews, ICF linking) or structural (e.g., factor analyses or Rasch), construct (known group differences; convergent/divergent associations), criterion (concurrent/predictive), responsiveness; predictive, evaluative or discriminative properties were established<br>3. Three or more indicators of reliability and validity were examined concurrently and provide a rich view on measurement properties. |
| | 1 | Two or more clinical measurement properties were evaluated, however, scope was narrow and did not meet above criteria. (e.g., internal consistency and one other indicator of validity or reliability ). |
| | 0 | The scope of clinical measurement properties was very narrow as indicated by a narrow evaluation of only one form of reliability or validity. |

| 5 | 2 | Authors performed a sample size calculation and obtained their recruitment targets. Post-doc power analyses and/or confidence intervals confirm that the sample size was sufficient to define relatively precise estimates of reliability or validity. |
|---|---|---|
|   | 1 | The authors provide an acceptable rationale for the number of subjects included in the study, but did not present specific sample size calculations or post-doc power analyses (or had a sample >100 but no justification). |
|   | 0 | Size of the sample was not rationalized or is clearly underpowered. |
| 6 | 2 | 90% or more of the patients enrolled for study were re-evaluated. |
|   | 1 | 70% or more of the enrolled patients were re-evaluated. |
|   | 0 | Less than 70% of the patients enrolled in the study were re-evaluated |
| **Measurements** | | |
| 7 | 2 | Documentation is provided for how the studied test is performed.  This includes adequate description of the measure/test and how it is administered or scored. The authors may provide or reference a published manual/article that outlines specific procedures for administration, scoring (including scoring algorithms, handling of missing data) and interpretation that included any necessary information about positioning/active participation of the client, any special equipment required, calibration of equipment if necessary, training required, cost, examiner procedures/actions. If no manual is available, then the text describes key details of procedures in sufficient detail so they could be replicated. |
|   | 1 | The test(s) and its administration procedures are referenced; but there is inadequate description of the test procedures. |
|   | 0 | Minimal description of test procedures without appropriate references. |

| 8 | 2 | This item addresses the overall study procedures for administering all study measures (study measure and its comparators) in an unbiased way. Test procedures should not introduce systematic errors in the estimation of the clinical measurement properties. This includes standardized procedures for who completed or administered the measures. For self-report, this includes order of presentation, who completed at what time interval; handling of missing items. If relevant, then the paper should include how cultural literacy issues were handled (e.g., exclusion, assisted or surrogate completion). For impairment measures, procedures would include calibration of any equipment; use of consistent measurement tools and scoring, a priori exclusion of any participants likely to give invalid results/unable to complete testing (not exclusion of after enrollment); use of standardized instructions and test procedures. This can include order of administration of test and quality checking of scores. For reliability testing, the appropriate retest interval will depend on the nature of the condition; but for acute conditions it may require retesting within 48 hours; whereas chronic/stable conditions are commonly retested within 4-14 days. For estimation of clinical change, retest intervals should be ones during which a meaningful clinical change would have occurred (and from an intervention with known effectiveness). The evaluator decides overall whether this has sufficiently been addressed by the methods described. |
|---|---|---|
| | 1 | No obvious sources of bias in the study test protocol or how tests were performed/administered is apparent; but there were suboptimal procedures or an inadequate description of the measurement protocol to be insured control of bias or that procedures were standardized. |
| | 0 | No description of the overall procedures for administering study tests; OR an obvious source of bias in data collection methods. |

**Analyses**

| 9 | 2 | Authors clearly defined which specific analyses were conducted for each of the stated specific hypotheses/questions of the study. This may be accomplished through organization of the results under specific subheadings or by demarcating which analyses addressed specific clinical measurement properties. Data was presented for each hypothesis/research question posed. |
|---|---|---|
| | 1 | Data was presented that addressed each of the measurement questions posed, but authors did not link specific analyses to specific research questions or hypotheses. |
| | 0 | Data was not presented for every hypothesis or clinical measurement property outlined in the purposes or methods. |

| 10 | 2 | <u>Tests selected</u> - Appropriate statistical tests were conducted to calculate a point estimate for clinical measurement properties.  Examples are provided below; but are not exhaustive.<br><br>1.  Reliability (Relative=ICCs (Shrout & Fleiss, 1979) for quantitative, Kappa (Landis & Koch, 1977) for nominal data); absolute (SEM or plot of score differences vs. average score showing mean and  2SD limit – as per Altman and Bland) (Bland & Altman, 1986; Bland & Altman, 1987)<br><br>2.  Clinical relevance - minimal detectable change, clinically important difference (Jaeschke, Singer, & Guyatt, 1989; Beaton et al., 2001; Wells et al., 2001)<br><br>3.  Validity<br><br>a. Validity associations - Pearson correlations for normally distributed data, Spearman rank correlations for ordinal data; or other correlations, if appropriate<br><br>b. Validity tests of significant difference - an appropriate global test like analysis of variance was used where indicated, with post-hoc tests that adjusted for multiple testing<br><br>c. Validity of items scaling/responses - Rasch analysis or item response (Baylor et al., 2011; Pallant & Tennant, 2007; Kyngdon, 2006; Cipriani, Fox, Khuder, & Boudreau, 2005; Smith, Jr., Conrad, Chang, & Piazza, 2002)<br><br>4. Responsiveness (Beaton, Bombardier, Katz, & Wright, 2001)- standardized response means or effect sizes or other recognized responsiveness indices were used. |
| | 1 | Appropriate statistical tests were used in some instances; but suboptimal choices were made in other analyses. |
| | 0 | Inappropriate use of statistical tests - incorrect tests for type of data; or a lack of analysis |
| 11 | 2 | The study goes beyond a single statistical point estimate of a clinical measurement property and providing supporting statistical analyses that increases confidence in the findings in terms of precision of the (key) indicator; or provide an alternate form of analysis of the clinical measurement property. The evaluator decides if these analyses are appropriate and informative.  For example, with reliability, at least 2 of the following would constitute appropriate and informative analysis beyond a point estimate a reliability coefficient: 1. confidence intervals around the point estimate; 2. Comparison to appropriate referenced benchmarks or standards; or 3. SEM or MDC.  For correlations, tests of significance or confidence intervals were presented and indicators of the criterion benchmarks were provided.  For studies involving cross-cultural validation, the analyses should compare multiple clinical measurement properties previously established for the measure and explain the extent to which the translated version is in accordance with these previously reported properties on the source measure. |

| | 1 | Either precision definition (confidence intervals) or appropriate benchmark comparison were used - NOT both. OR Some analyses were associated with indicators of precision or alternate form of analysis -but not all key indicators. |
|---|---|---|
| | 0 | Inappropriate use of benchmarks or confidence intervals; or indicators of precision or alternate form are absent |
| **Recommendations** | | |
| 12 | 2 | Authors made specific conclusions and clinical measurement recommendations that were clearly related to each hypotheses/question posed in the study and that were supported by the data presented.  Ideal recommendations would state the estimated status of the clinical measurement property, the confidence in the estimate and the context for which those apply.  To achieve a 2, the conclusion must be specific; and conclusions cannot overstate the clinical measurement properties observed the study; nor ignore suboptimal measurement properties found. |
| | 1 | Authors made conclusions and clinical measurement recommendations that were basically true (supported by study data); but vague. That is, they do not specify the extent, confidence or context of the findings.  (The measure is "reliable and valid ") OR authors made specific clinical measurement recommendations; but for only some of the study hypotheses. |
| | 0 | Authors did not make conclusions about clinical measurement; OR made recommendations that were in contradiction to the actual data presented |

© MacDermid 2011

**List with excluded studies with reasons**

| | |
|---|---|
| 1. Abbott et al 2014 | Ineligible population |
| 2. Beattie et al 2011 | Ineligible population (less than 50%) |
| 3. Hoeskstra et al 2014 | No properties for GRoC scales |
| 4. Chansirinukor 2019 | No properties for GRoC scales |
| 5. Chien et al 2015 | No properties for GRoC scales |
| 6. Cruz et al. 2015 | No properties for GRoC scales |
| 7. Foroutani et al 2018 | No English (Persian language) |
| 8. Gagnon et al 2018 | Ineligible population |
| 9. Hefford et al 2012 | Ineligible population |
| 10. Hung et al 2019 | Ineligible population |
| 11. Sharma et al 2017 | Ineligible population |
| 12. Stevens et al 2019 | Ineligible population |
| 13. Meyer et al 2014 | Ineligible population |

# 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 |
| **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. | 3-5 |
| Objectives | 4 | Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS). | 4-5 |
| **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. | 5 |
| 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 |
| 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. | 6 |
| Search | 8 | Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated. | Appendix1 |
| 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). | 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. | 6-7 |
| Summary measures | 13 | State the principal summary measures (e.g., risk ratio, difference in means). | 8-9 |
| 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. | 8-9 |

# 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). | 8-9 |
| Additional analyses | 16 | Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified. | 8=9 |
| **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. | 9 |
| 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. | 9-10 |
| 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 |
| 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. | 10-12 |
| Synthesis of results | 21 | Present results of each meta-analysis done, including confidence intervals and measures of consistency. | 13 |
| Risk of bias across studies | 22 | Present results of any assessment of risk of bias across studies (see Item 15). | 10 |
| Additional analysis | 23 | Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]). | 13 |
| **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). | 14-15 |
| 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 |
| Conclusions | 26 | Provide a general interpretation of the results in the context of other evidence, and implications for future research. | 14-15 |
| **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. | 18 |

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

Page 2 of 2

# Psychometric Properties of the Global Rating of Change Scales in Patients with Neck Disorders: A Systematic Review with Meta-Analysis and Meta-Regression

**SCHOLARONE™**
Manuscripts

1  **Psychometric Properties of the Global Rating of Change Scales in Patients with Neck**

2  **Disorders: A Systematic Review with Meta-Analysis and Meta-Regression**

3  Pavlos Bobos[1], Joy C MacDermid[2], Goris Nazari[3], Rochelle Furtado[4] and CATWAD co-authors[5]

4

5  [1]Pavlos Bobos PT, PhD(c), (corresponding author) Doctoral Candidate, Western's Bone and Joint

6  Institute, Department of Health and Rehabilitation Sciences, Western University, Elborn College,

7  1201 Western Road, N6G 1H1, London, Ontario, Dalla Lana School of Public Health, Institute of

8  Health Policy Management and Evaluation, Department of Clinical Epidemiology and Health Care

9  Research, University of Toronto, Canada, (pbobos@uwo.ca), tel: +1 519 661 2111 x88912

10  [2]Joy C MacDermid BScPT, PhD,  Professor, Physical Therapy and Surgery, Western University,

11  London, ON and Co-director Clinical Research Lab, Hand and Upper Limb Centre, St. Joseph's

12  Health Centre, London, Ontario; Professor Rehabilitation Science McMaster University,

13  Hamilton, ON, Canada (jmacderm@uwo.ca)

14  [3]Goris Nazari PT, PhD(c) Doctoral Candidate, Western's Bone and Joint Institute, School of

15  Physical Therapy, Department of Health and Rehabilitation Sciences, Western University,

16  London, Ontario, Canada, (gnazari@uwo.ca)

17  [4]Rochelle Furtado MSc Western's Bone and Joint Institute, School of Physical Therapy,

18  Department of Health and Rehabilitation Sciences, Western University, London, Ontario, Canada,

19  (rfurtad5@uwo.ca)

20  [5]**CATWAD:** Michele Sterling m.sterling@uq.edu.au, Anne Söderlund anne.soderlund@mdh.se,

21  Michele Curatolo, curatolo@uw.edu, James M Elliott j-elliott@northwestern.edu, David  Walton

22  dwalton5@uwo.ca, Helge Kasch helgkasc@rm.dk, Linda Carroll linda.carroll@ualberta.ca,

23  Hans Westergren Hans.Westergren@skane.se, Gwendolen Jull g.jull@uq.edu.au, Eva-Maj

24  Malmström eva-maj.malmstrom@med.lu.se, Luke B Connelly l.connelly@uq.edu.au, Joy C

25  MacDermid jmacderm@uwo.ca, Mandy Nielsen mandy.nielsen@griffith.edu.au, Pierre Côté

26  pierre.cote@uoit.ca, Tonny Elmose Andersen tandersen@health.sdu.dk, Trudy Rebbeck

27  trudy.rebbeck@sydney.edu.au, Annick Maujean a.maujean@uq.edu.au, Sarah Robins

28  s.robins1@uq.edu.au, Kenneth Chen k.chen8@uq.edu.au, Julia Treleaven j.treleaven@uq.edu.au

30  **Word count:** 3908

1

31  **ABSTRACT**

32  **Objective:** The purpose of this systematic review was to critically appraise and synthesize the

33  psychometric properties of Global Rating of Change (GROC) scales for assessment of patients

34  with neck pain.

35  **Design:** Systematic review

36  **Data sources:** A search was performed in 4 databases (MEDLINE, EMBASE, CINAHL,

37  SCOPUS) until February 2019.

38  **Data extraction and synthesis:** Eligible articles were appraised using Consensus-based Standards

39  for the selection of health Measurement Instruments (COSMIN) checklist and the Quality

40  Appraisal for Clinical Measurement Research Reports Evaluation Form.

41  **Results:** The search obtained 16 eligible studies and included in total 1533 patients with neck pain.

42  Test-retest reliability of Global Perceived Effect (GPE) was very high (Intra-class correlation

43  coefficient (ICC) = 0.80 to 0.92) for patients with whiplash. Pooled data of Pearson's r indicated

44  that GROC scores were moderately correlated with neck disability change scores (0.53, 95% CI:

45  0.47 to 0.59). Pooled data of Spearman's correlations indicated that GROC scores were moderately

46  correlated with neck disability change scores (0.56, 95% CI: 0.41 to 0.68).

47  **Conclusions:** This study found excellent quality evidence of very good to excellent test-retest

48  reliability of GPE for patients with Whiplash Associated Disorders. Evidence from very good-to-

49  excellent quality studies found that GROC scores are moderately correlated to an external criterion

50  patient-reported outcome (PROM) measure evaluated pre-post treatment in patients with neck

51  pain. No studies were found that addressed the optimal form of GROC scales for patients with

52  neck disorders or compared the GROC to other options for single-item global assessment.

53  **Prospero registration number:** CRD 42018117874

54

2

**Strengths and limitations of this study**

- We rated the quality of individual studies and the overall risk of bias using two standardized approaches

- Our focus on neck pain increased the specificity of results but are not necessarily applicable to other musculoskeletal conditions

- Conceptual concerns about global ratings of change being affected by recall bias are not adequately addressed by psychometric evidence

- No studies addressing the optimal form of global rating were found.

**Introduction**

Neck pain is the 4th leading cause of disability and approximately half of adult the population with neck pain will experience a clinically important episode once in their lifetime. [1–3] The annual prevalence of neck pain it is estimated between 15% and 50%, with females having a higher prevalence rate than males. [2,3] Neck pain has been associated with many other comorbidities such as headaches, dizziness, anxiety, depression, back pain and arthralgias.[3–6] Several different methods for classifying neck pain have been described, using indicators such as duration (acute, sub-acute or chronic), degree of interference (low, moderate, severe) or most likely structure at fault (e.g. neuropathy vs. mechanical). [7]

As part of a patient-centric approach to care, clinicians will commonly evaluate response to intervention by asking the patient directly whether they feel better, worse, or the same since the prior encounter. While direct questioning can provide a qualitative indicator of change in status, many best practice guidelines endorse use of some form of quantified patient-reported outcome (PRO) as an adjunct to oral self-report. PROs are available to quantify several different constructs

3

79 in people with neck pain, including pain severity, disability and neck function. [8] Any PRO

80 intended to provide an estimate of change over time should be responsive to subtle shifts in the

81 patient's condition. To facilitate interpretation of change scores, a common property of many such

82 tools is the minimum clinically important difference (MCID), which is a change threshold that

83 corresponds to the minimum shift in scale values that most patients would indicate corresponds to

84 an important change in their overall condition. A well-recognized approach to establishing an

85 MCID for a PRO is to compare the magnitude of change against an anchor, most commonly a

86 Global Rating of Change (GROC) scale. These scales allow patients or study participants to

87 indicate whether their condition has gotten worse, better, or stayed the same and to quantify the

88 magnitude of that change. As they have been adopted as a sort of 'standard' against which change

89 in other tools is compared, the GROC can also be used on its own as an omnibus generic indicator

90 of change. [8]

91 Despite being accepted as a standard measure, there is considerable variation in how the

92 GROC has been constructed and implemented in research in neck pain. GROC scales consist of

93 ordered categories which may have different ranked levels (some have 15 levels, some 11 levels,

94 and others have 7 levels). The common structure across these is the use of a middle '0' score

95 corresponding to 'no change', with negative values indicating magnitudes of worsening while

96 positive values indicate improvement.[9] Variations of the GROC (in name or structure) include

97 the "Global Perceived Effect", "Patient Global Impression of Change", "Transition Ratings", and

98 "Global Scale". [9]

99 A well-established component of health outcomes is having a tool with strong

100 psychometric properties of validity, reliability and responsiveness to be able to monitor change.

101 While recent research [8] has examined the psychometric properties of the most commonly

4

102 reported PROs for neck disorders, to date there has been no systematic review to summarize the

103 measurement properties of GROC scales themselves in patients with neck disorders. Therefore,

104 this systematic review aims to critically appraise and synthesize the psychometric properties of the

105 GROC scales in patients with neck disorders.

106

107 **METHODS**

108 *Patient and Public Involvement*

109 There was no patient or public involvement in the design or planning of this study.

110

111 *Study Design and Protocol Registration*

112 We conducted a systematic review to evaluate the psychometric properties of GROC scales in

113 patients with neck disorders. The protocol was registered in PROSPERO register database with

114 registration number: CRD 42018117874

115

116 *Eligibility Criteria*

117 We included studies in this systematic review if the following criteria were met [10–12]:

- 118  • Design: psychometric testing, randomized/ cohort studies

- 119  • Participants: > 50% of the study's patient population with neck conditions/disorders,

- 120  • Intervention/Comparison: studies that reported on the psychometric properties (reliability,

121  validity, responsiveness) of GROC, Global Perceived Effect (GPE) and Patient Global

122  Impression of Change (PGIC),

- 123  • Outcomes: GROC, GPE and PGIC

- 124  • Articles were written in English language only

5

125 Studies with no data on the GROC scales' psychometric properties, and conference

126 abstract/posters were excluded from this systematic review.

127

*Information Sources*

129 To identify studies on the psychometric properties (reliability, validity, responsiveness) of the

130 GROC, GPE and PGIC we searched the Medline, EMBASE, Scopus and CINAHL databases from

131 inception till February 2019, using a combination of keywords. Furthermore, we identified

132 additional studies by examining the reference list of each of the selected studies. The full list with

133 keyword strategy is presented in **APPENDIX 1**.

134

*Study Selection*

136 Two investigators (PB and GN) performed the systematic electronic searches independently in

137 each database. The same investigators then proceeded to identify and remove the duplicate studies.

138 In the next stage, we performed the independent screening of the titles and abstracts and any full-

139 text article marked as include or uncertain were obtained. In the final stage, the same two

140 independent authors performed the full text reviews independently to assess final article eligibility.

141 In case of disagreement, a third reviewer; the most experienced member (JM), facilitated a

142 consensus through discussion.

143

*Data Extraction*

145 The fourth author (RF) performed the data extractions. The extracted data were then cross-checked

146 by another author (PB). Data extraction included the author, year, study population/condition,

147 setting, sample size, age, properties evaluated, retest-interval, and the intervention protocol (if used

6

148   to assess responsiveness parameters). [13,14] For reliability estimates, Standard Error of

149   Measurement (SEM), Intra-class Correlation Coefficient (ICC), Minimal Detectable Change

150   (MDC) and 95% confidence intervals were extracted. [13,14] The ICC interpretation of ICC < 0.40

151   indicating poor, $0.40 \leq ICC < 0.75$ indicating fair-to-good and $ICC \geq 0.75$ indicating excellent

152   reliability were used as a common benchmark.[15] For validity estimates, correlation coefficient

153   (Pearson's/Spearman) and the 95% confidence intervals were extracted. [13,14] Evan's guidelines

154   to interpret the strength of the correlation was used which included: 0.00–0.19 "very weak", 0.20–

155   0.39 "weak", 0.40–0.59 "moderate", 0.60–0.79 "strong", and 0.80–1.00 "very strong". [16] For

156   responsiveness estimates, the Effect Size (ES), Standardized Response Mean (SRM), Clinically

157   Important Difference (CID), and/or Minimal Clinically Important Difference (MCID) including

158   the method of MCID estimation – Anchor-/Distribution-based methods, and 95% confidence

159   intervals were extracted. [13,14] To assist clinical decision making, standard benchmark scores of

160   trivial (< 0.20), small ($\geq 0.20$ to < 0.50), moderate ($\geq 0.50$ to < 0.80) or large ($\geq 0.80$), as proposed

161   by Cohen, were used. [17] When insufficient data were presented, PB contacted the authors by

162   email and requested further data.

163

164   *Consensus-based Standards for the selection of health Measurement Instruments (COSMIN)*

165   Consensus-based Standards for the selection of health Measurement Instruments (COSMIN)

166   assesses the risk of bias for the psychometric properties reported on a property-by-property basis.

167   A score for the risk of bias in estimates of psychometric properties was assessed by two authors

168   (PB) and (RF) using the new (COSMIN) checklist.[18] If disagreement was present a third person

169   (JM) assist in resolving the discrepancy. Each study was assessed by COSMIN on the 4-point scale

170   as "very good", "adequate", "doubtful" or "inadequate" for each of the checklist criteria for

7

171 relevant measurement properties (e.g. reliability, responsiveness, etc.). According to COSMIN,

172 when determining the overall score for each measurement property, the worst score counts method

173 was used wherein the lowest score for the checklist criteria of the relevant property was taken as

174 the overall score. [19] We then assessed the result of individual studies on a measurement property

175 against the updated criteria for good measurement properties. This involved the evaluation of

176 results of included studies as either sufficient (+), insufficient (−), or indeterminate (?). [18]

177

178 *Quality Appraisal for Clinical Measurement Research Reports Evaluation Form*

179 A summary score for the overall quality of individual studies was appraised independently by the

180 authors (PB) and (RF) using a structured clinical measurement specific appraisal tool. [13,14] In

181 case of disagreement a third person was consulted (JM) to resolve the conflict. The evaluation

182 criteria of this tool included twelve items: 1) Thorough literature review to define the research

183 question; 2) Specific inclusion/exclusion criteria; 3) Specific hypotheses; 4) Appropriate scope of

184 psychometric properties; 5) Sample size; 6) Follow-up; 7) The authors referenced specific

185 procedures for administration, scoring, and interpretation of procedures; 8) Measurement

186 techniques were standardized; 9) Data were presented for each hypothesis; 10) Appropriate

187 statistics-point estimates; 11) Appropriate statistical error estimates; and 12) Valid conclusions

188 and recommendations. [13,14] An article's total score − quality - was calculated by the sum of

189 scores for each item, divided by the numbers of items and multiplied by 100%. [13,14] Overall,

190 the quality summary of appraised articles range from (0%-30%) Poor, (31%-50%) Fair, (51%-

191 70%) Good, (71%-90%) Very Good, and (>90%) Excellent. [13,14]

192

193 *Synthesis of Results*

194 A qualitative synthesis was conducted to report findings on test-retest reliability statistics. A meta-

195 analysis of Pearson's and Spearman's correlation was performed in R (version 3.6.1) with

196 metaphor package.[20] The meta-analyses were conducted using a random effect model and the

197 correlation coefficients were converted to z values. Heterogeneity was deemed substantial if $I^2$

198 values were more than 50%. [21] A Meta-regression was planned to explore the sources of

199 unexplained heterogeneity by considering the following factors: a. neck pain with or without

200 radicular symptoms, b. acute or chronic, c. age and d. sex. Forest plots were created using means

201 and 95% confidence intervals for correlation coefficients. We summarize the main results of the

202 included articles based on the neck disorders, reported psychometric estimate and the study quality

203 ratings.

204

205 **RESULTS**

206 *Study Selection*

207 Our search yielded 8,837 articles. After removal of duplicates, 6,027 studies remained and were

208 screened using their title and abstract; leaving 29 articles selected for full-text review. Of these, 16

209 studies were considered eligible. [22,23,24–31,32–37] The flow of the study selection process is

210 presented in **Figure 1.**

211

212 *Study Characteristics*

213 The 16 eligible studies were conducted between 2006 and 2017 and included 1533 participants

214 with neck pain/disorders (mean of 96 participants per study). [22,23,24–31,32,34–37,] Study size

215 ranged from 29 to 200 participants. A summary description of all the studies included is displayed

216 in **Table 1.** Concurrent validity was evaluated in 14 studies by comparing the difference of pain

217  intensity, disability and function scores with the score of GROC scales. Two studies [26,31]

218  examined the test-retest reliability of a 7-point and an 11-point GPE scale for patients with

219  whiplash-associated disorders (WAD). One study [24] examined whether occurrences of within-

220  and between-session changes were significantly associated with functional outcomes, pain, and

221  self-report of recovery in patients at discharge who were treated with manual therapy for

222  mechanical neck pain.

223

224  *COSMIN Risk of Bias rating and Quality appraisal of the Included Studies*

225  Regarding the risk of bias, all studies were rated as very good (**Table 2**). The quality of the studies

226  ranged from 88% to 96% (**Table 3**). The most common flaws were 1) lack of/inadequate sample

227  size calculations, 2) missing data (i.e. inadequate follow up), and 3) inconsistencies between the

228  data presented and hypothesis stated.

229

230  *Reported GROC scales*

231  The most commonly reported GROC scale (n=6 studies) was a 15-point scale with the most

232  frequent anchors being "-7 (a very great deal worse) to zero (about the same) to +7 (a very great

233  deal better)". A 7-point scale was reported in 5 studies, 11- and 5-point scales were reported in 2

234  studies and a 9-point scale in one study. The anchors in those scales varied greatly and are

235  presented in Table 1. Only 6 studies [26,31–33,35,36] reported full detail regarding the specific

236  questions asked of the patients with neck disorder when a GROC scale was administered. Those

237  questions that were reported are presented in **Box 1.**

238

239

10

*Reliability Measures*

240

241 Two studies were included that examined test-retest reliability of GPE for patients with WAD.

242 Kamper et al. (2010) [26]  examined the [time interval] test-retest reliability of an 11-point GPE

243 scale in 134 patients with chronic WAD and reported an Intra-class Correlation Coefficient (ICC)

244 of 0.99 (95% CI 0.99 to 0.99) at baseline, 0.96 (0.95 to 0.97) at 6 weeks, and 0.92 (0.89 to 0.94)

245 at 12 months (**Table 4**). Ngo et al. (2010) assessed the test-retest reliability of a 7-point scale of

246 GPE in patients with acute WAD at 3 to 5 days. [31] The ICC and 95% confidence intervals (CI)

247 were used to determine the test–retest reliability of the two versions of the perceived recovery

248 questions using their original seven-item responses. Ngo et al. also computed weighted kappa

249 coefficients and 95% CI using quadratic weights to determine whether the distribution of responses

250 influenced the reliability as measured by the ICC. An ICC for general recovery of 0.70 (0.60 to

251 0.80) and an ICC for neck pain questions of 0.80 (0.72 to 0.87) were found. A weighted Kappa

252 was also calculated (Kappa = 0.70 (0.42 to 0.98)) at six weeks for general recovery and at six

253 weeks Kappa = 0.80 (0.51 to 1.0) for neck pain questions (**Table 4**).

254

*Validity Measures*

255

256 We found 14 studies that examined concurrent validity measures between GROC and another

257 PRO.[22,23,25,27–30,32,34,35,36–38] Correlations of Pearson's and Spearman's coefficients

258 between GROC and another PRO were ranging from very weak to very strong correlations. The

259 validity measures are presented and summarized in Table 5.

260

261

262 *Meta-Analysis and Meta-Regression of Correlations between Disability change scores and GROC*

263 *scores*

264 Five studies [23,25,34,37,38] of very good-to-excellent quality reported the Pearson correlation

265 coefficients between neck disability change scores and the GROC scores and were pooled together.

266 We found that GROC was positively correlated with disability change scores (r = 0.53, 95% CI:

267 0.47 to 0.59, $I^2$ = 0%). Six studies [27–30,32,36] of very good-to-excellent quality reported the

268 Spearman correlation coefficients between neck disability changes scores and the GROC scores

269 and were pooled together. We found that GROC was moderately correlated with disability change

270 scores (rho = 0.56, 95% CI: 0.41 to 0.68, $I^2$= 85%). The forest plots with correlation coefficients

271 with 95% CIs are presented in Figure 2-3. Our meta-regression showed that age was found as a

272 significant factor in influencing Fisher's Z scores (β = -0.034, 95% CI -0.05 to -0.01, p = 0.001).

273 The model explained 68% of the variance ($R^2$ = 0.68) (Figure 4).

274

275 *Area under the curve (AUC) – Sensitivity and Specificity*

276 Cook et al. [24] found that between-session NPRS- pain changes were associated with greater than

277 3-point change on the GROC at 96-hours (AUC=0.76). The pain change associated with GROC

278 was more specific (Specificity=79.2%, range: 62.2 - 91.1) than sensitive (Sensitivity=65.6%,

279 range: 57.9 to 74.6). Those with a 36.7% between-sessions change in pain were also 7.3 times

280 more likely to report an improvement of greater than 3 points change on the GROC than those

281 who did not achieve a 36.7% change in pain (**Table 4**).

282

283 **DISCUSSION**

284    This review has synthesized the current research from 16 studies that aimed to evaluate the

285    psychometric properties of GROC scales for patients with neck disorders, with the goal to provide

286    evidence for clinicians and researchers concerning its use within clinical practice and research.

287    From the 16 included studies, only 2 studies [26,31] reported test-retest reliability statistics of the

288    7- and 11-ranked categories of GPE scales for patients with WAD only. We were able to pool data

289    from 12 studies regarding concurrent validity of GROC scales and neck disability change scores

290    at one time point after the interventions. Themes influencing interpretation of the GROC were

291    explored in a study [33] that evaluated the factors that contribute to how patients respond to a

292    question on global perceived effect. This study found that treatment process, biomechanical

293    performance, self-efficacy and the nature of the condition may influence the responses on global

294    perceived effect, which is consistent with what we would expect for patients with neck pain. This

295    suggests that change is a complex multifactorial global concept. A strength of GROC is that it is

296    intended as a global assessment, and it can be assumed that it reflects the aspects of change

297    important to the individual patient.

298    Reliability can be defined as the degree to which a measure produces consecutive results

299    with the least amount of random error when the status of the population remains unchanged. The

300    reliability of GPE displayed an excellent test-retest reliability of ICC>0.90 over an interval of 6

301    weeks and 12 months for patients with WAD. Conducting an assessment with a long test-retest

302    interval (e.g. 12 months), can provide challenges as there is higher risk of individuals with WAD

303    being symptomatically unstable.[9] Determining if patients are symptomatically-stable can be

304    achieved by administering another PRO such as the Single Assessment Numeric Evaluation

305    (SANE)[39], however, the 7- and 11- ranked categories of GPE scales still demonstrated good

306    stability properties at long test intervals (i.e., of 6 weeks and 12 months).[26] Therefore, the

13

307 measurements of the reliability parameters of the GPE may be very useful during longer test

308 intervals in clinical trials.

309       The psychometric property of validity is defined as the degree to which a PRO measures

310 what it is intended to measure. Pooled data from 11 studies overall suggest that post-treatment

311 changes of on validated disability outcome measures were moderately (Pearson's r = 0.51, 95%

312 CI: 0.43 to 0.58; Spearman's rho = 0.56, 95% CI: 0.41 to 0.68) correlated to change in perceived

313 effect) (Figure 2-3). This finding suggests that GROC scores taken at one point in time were related

314 to scores in pain and disability in patients with neck disorders, as measured by standardized

315 measures taken at 2 points in time. We identified one study [24] that found a 36.7% change in pain

316 for within- and between- session changes was associated with a 50% reduction in the NDI and an

317 improvement of >3 levels on a 15-ordinal level GROC scale for patients with neck pain. This

318 quantified predictive change value may have clinical utility for use in clinical practice.

319       Previous studies [9,40] have indicated serious concerns about the conceptual validity of the

320 global rating of change. The review by Kamper et al.[9]  clearly showed that GROC was related

321 to final status more than change and was least related to baseline health status. This result

322 undermines the premise of what the global rating of change actually measures. For this reason, we

323 conclude that the 0.50 pooled correlation across 12 studies between the GROC and other PROM

324 change scores (e.g. Neck Disability Index (NDI) scores) may reflect a relationship between follow-

325 up status and change rather than supporting the contention that GROC actually measures change.

326 This would also explain why only 25% of the variation in GROC change scores was explained by

327 changes scores from a PROM change score measured at 2 points in time. In all studies, participants

328 completed the GROC scale at one time point after the intervention, and hence recall bias is a cause

329 for concern. However, another potential factor for moderate correlations is that the PROMs that

14

330    have been used as the comparator with GROC scores may not reflect priorities that are important

331    to patients. That is, the field has largely been driven by assumptions that the GROC is a 'gold

332    standard' for evaluating true change in a respondent's condition or status, and that all items on the

333    comparator PROM are of equal importance to all people with that condition. The work presented

334    herein challenges the valorization of the GROC as a gold standard for change, and prior work has

335    challenged the notions that all PROM items are equally important.[9,41,42] It is therefore possible

336    that the very constructs being evaluated require greater critical discourse before authors can say,

337    with confidence, that one scale functions well or poorly based on its associations with another

338    scale. Since no studies compared a retrospective global assessment of the GROC to pre-post single

339    item global PROM e.g. the SANE, we do not know the extent to which these two factors

340    contributed to moderate correlation.

341        A unique aspect of this study was that it focused on global rating of change scales in a neck

342    pain patient population. Our study appraisal suggests that future studies concerning GROC should

343    include adequate sample sizes, maintain a rigorous follow up and report appropriate statistical error

344    estimates, since these were often inadequate. Various critical appraisal tools exist, and the

345    perspectives and ratings may differ across instruments. COSMIN is just one methodology that can

346    be used to synthesize or evaluate outcome measures and other methods might be equally valid or

347    provide different perspectives. We used 2 different critical appraisal tools to evaluate quality from

348    2 perspectives. The COSMIN risk of bias assessments reflects the level of confidence in the

349    conclusions and pooled estimates. The quality appraisal tool focuses on design issues in the studies

350    and reflects gaps in research designs that should be considered in interpretation of current research

351    and improved in future studies. Substantial heterogeneity was detected ($I^2 > 50\%$) in pooled

352    Spearman's correlation coefficients which is a concern when pooling data. Sources of the observed

15

353 heterogeneity were identified in our meta-regression results. Our univariate meta-regression

354 analysis indicated that age across the studies explained 68% of the variance (**Figure 4**). Other

355 factors such as type of neck pain (with or without radicular symptoms), acute or chronic and sex

356 did not explain the remaining heterogeneity (not statically significant). In our meta-regression, we

357 used a patient level characteristic to identify the observed heterogeneity and therefore, our model

358 may be vulnerable to aggregation bias. Furthermore, the scope of our literature search was focused

359 on identifying full-text papers written in English only.

360 While this study included 16 studies, only 2 of these reported reliability statistics for GROC

361 scales for patients with chronic WAD. Therefore, the applicability of our study is mostly limited

362 to patients with chronic WAD. For validity measurements, GROC scales were mostly investigated

363 by correlation analyses to evaluate the external responsiveness of another PRO measure over a

364 specific time point. From our meta-analysis, we can be confident that the GROC scores were

365 moderately correlated with neck disability change scores. However, more robust psychometric

366 design studies to test the measurement properties of GROC scales as the primary outcome of

367 investigation are highly needed. Future studies should aim to test to what extent the different range

368 of items (e.g. 7-level scale vs 11-level scale), the anchors (e.g. much worse vs much better) may

369 affect the measurement properties of GROC scales for patients with neck disorders. Also, it is

370 important to indicate that most outcome measures are ordinal and assume that additive scores of ordinal

371 items can be treated as interval level. This potentially could lead to scaling problems even in the face of

372 strong psychometric properties. The main protection we have is to create new scales or retrofit existing

373 scales based on Rasch analysis. Also, we acknowledge that the majority of work done on the GROC scales

374 has been performed using statistical approaches that are most appropriate to linear rather than ordinal data

375

376 **CONCLUSIONS**

16

377 This study found excellent quality evidence of very good to excellent test-retest reliability of GPE

378 for patients with WAD. Evidence of very good to excellent quality studies found that GROC scores

379 are moderately correlated to an external criterion PROM measure measured pre-post treatment in

380 patients with neck disorders. Studies addressing the optimal form of GROC scales for patients with

381 neck disorders or comparing the GROC to other options for single-item global assessment of

382 change were not found.

383

384 **Authors' contributions**

385 PB contributed significantly to conception and design of the study, data extraction, critical

386 appraisal, interpretation of data and drafting of the manuscript. GN, and RF were involved in

387 literature search, critical appraisal and interpretation of data and drafting. GN was involved in

388 critical appraisal and drafting. JM was also involved in the conception and design of the study,

389 drafting, and revised the manuscript for important intellectual content. JM and CATWAD were

390 involved in the drafting and review of the manuscript. All authors have given their final approval

391 on the manuscript to be published

392

393 **Declarations**

394 **Ethics approval and consent to participate**

395 Not applicable

396 **Consent for publication**

397 Not applicable

398 **Availability of data and material**

399 Data sharing is not applicable to this article as no datasets were generated or analyzed during the

17

400  current study

**Funding Statement**

**Competing Interest Statement**

405  None to report

406

407

**References**

409  1  Murray CJL, Abraham J, Ali MK, *et al.* The State of US health, 1990-2010: Burden of diseases,

410  injuries, and risk factors. *JAMA - J Am Med Assoc* Published Online First: 2013.

411  doi:10.1001/jama.2013.13805

412  2  Fejer R, Kyvik KO, Hartvigsen J. The prevalence of neck pain in the world population: A

413  systematic critical review of the literature. Eur. Spine J. 2006. doi:10.1007/s00586-004-0864-4

414  3  Hogg-Johnson S, van der Velde G, Carroll LJ, *et al.* The Burden and Determinants of Neck Pain in

415  the General Population. Results of the Bone and Joint Decade 2000-2010 Task Force on Neck

416  Pain and Its Associated Disorders. *J Manipulative Physiol Ther* Published Online First: 2009.

417  doi:10.1016/j.jmpt.2008.11.010

418  4  Bobos P, Nazari G, Palimeris S, *et al.* The contribution of health and psychological factors in

419  patients with chronic neck pain and disability: A cross-sectional study. *J Clin Diagnostic Res*

420  2018;**12**:YC04–7. doi:10.7860/JCDR/2018/31284.11203

421  5  MacDermid JC, Walton DM, Bobos P, *et al.* A Qualitative Description of Chronic Neck Pain has

422  Implications for Outcome Assessment and Classification. *Open Orthop J* 2017;**10**:746–56.

423  doi:10.2174/1874325001610010746

424  6  Treleaven J. Sensorimotor disturbances in neck disorders affecting postural stability, head and eye

18

425    movement control-Part 2: Case studies. *Man Ther* 2008;**13**:266–75.

426    doi:10.1016/j.math.2007.11.002

427  7  Cohen SP. Epidemiology, diagnosis, and treatment of neck pain. *Mayo Clin Proc* 2015;**90**:284–99.

428    doi:10.1016/j.mayocp.2014.09.008

429  8  Bobos P, Macdermid JC, Walton DM, *et al.* Patient-reported outcome measures used for neck

430    disorders: An overview of systematic reviews. J. Orthop. Sports Phys. Ther. 2018;**48**:775–88.

431    doi:10.2519/jospt.2018.8131

432  9  Kamper SJ, Maher CG, Mackay G. Global Rating of Change Scales: A Review of Strengths and

433    Weaknesses and Considerations for Design. *J Man Manip Ther* 2009;**17**:163–70.

434    doi:10.1002/mus.21062

435  10  Nazari G, Bobos P, MacDermid JC, *et al.* The Effectiveness of Instrument-Assisted Soft Tissue

436    Mobilization in Athletes, Participants Without Extremity or Spinal Conditions, and Individuals

437    with Upper Extremity, Lower Extremity, and Spinal Conditions: A Systematic Review. *Arch Phys*

438    *Med Rehabil* Published Online First: February 2019. doi:10.1016/j.apmr.2019.01.017

439  11  Bobos P, Nazari G, Szekeres M, *et al.* The effectiveness of joint-protection programs on pain,

440    hand function, and grip strength levels in patients with hand arthritis: A systematic review and

441    meta-analysis. *J Hand Ther* 2018;**32**:194–211. doi:10.1016/j.jht.2018.09.012

442  12  Nazari G, Bobos P, MacDermid JC, *et al.* Psychometric properties of the Zephyr bioharness

443    device: A systematic review. *BMC Sports Sci Med Rehabil* 2018;**10**. doi:10.1186/s13102-018-

444    0094-4

445  13  Law MC, MacDermid J. *Evidence-based rehabilitation : a guide to practice.* Thorofare, NJ: :

446    Slack Incorporated 2014.

447  14  Roy JS, Desmeules F, MacDermid JC. Psychometric properties of presenteeism scales for

448    musculoskeletal disorders: A systematic review. *J Rehabil Med* Published Online First: 2011.

449    doi:10.2340/16501977-0643

450  15  Sańchez J. Rosner, B.: Fundamentals of Biostatistics, third edition. PWS-Kent, Boston 1990, xv,

19

451    655 pp., ISBN 0-534-91973-1. *Biometrical J* 1993;**35**:150. doi:10.1002/bimj.4710350205

452    16    Wuensch KL, Evans JD. Straightforward Statistics for the Behavioral Sciences. *J Am Stat Assoc*

453          Published Online First: 2006. doi:10.2307/2291607

454    17    Cohen J. Statistical power analysis for the behavioral sciences. Stat. Power Anal. Behav. Sci.

455          1988. doi:10.1234/12345678

456    18    Mokkink LB, de Vet HCW, Prinsen CAC, *et al.* COSMIN Risk of Bias checklist for systematic

457          reviews of Patient-Reported Outcome Measures. *Qual Life Res* Published Online First: 2018.

458          doi:10.1007/s11136-017-1765-4

459    19    Terwee CB, Mokkink LB, Knol DL, *et al.* Rating the methodological quality in systematic reviews

460          of studies on measurement properties : a scoring system for the COSMIN checklist. 2012;:651–7.

461          doi:10.1007/s11136-011-9960-1

462    20    Viechtbauer W. Conducting meta-analisys in R with metafor package. *J Stat Softw* 2010;**36**:1–48.

463    21    Higgins JPT, Thompson SG, Deeks JJ, *et al.* Measuring inconsistency in meta-analyses. *BMJ*

464          Published Online First: 2003. doi:10.1136/bmj.327.7414.557

465    22    Cleland J, Fritz J, Whitman J, *et al.* The reliability and construct validity of the Neck Disability

466          Index and Patient Specific Functional Scale. *Spine (Phila Pa 1976)* 2006;**31**:598–602.

467    23    Cleland JA, Childs JD, Whitman JM. Psychometric Properties of the Neck Disability Index and

468          Numeric Pain Rating Scale in Patients With Mechanical Neck Pain. *Arch Phys Med Rehabil*

469          2008;**89**:69–74. doi:10.1016/j.apmr.2007.08.126

470    24    Cook C, Lawrence J, Michalak K, *et al.* Is there preliminary value to a within- and/or between-

471          session change for determining short-term outcomes of manual therapy on mechanical neck pain?

472          *J Man Manip Ther* 2014;**22**:173–80. doi:10.1179/2042618614y.0000000071

473    25    Guzy G, Vernon H, Polczyk R, *et al.* Psychometric validation of the authorized Polish version of

474          the Neck Disability Index. *Disabil Rehabil* 2013;**35**:2132–7. doi:10.3109/09638288.2013.771706

475    26    Kamper SJ, Ostelo RWJG, Knol DL, *et al.* Global Perceived Effect scales provided reliable

476          assessments of health transition in people with musculoskeletal disorders, but ratings are strongly

477     influenced by current status. *J Clin Epidemiol* 2010;**63**:760-766.e1.

478     doi:10.1016/j.jclinepi.2009.09.009

479   27   Shaheen AAM, Omar MTA, Vernon H. Cross-cultural adaptation, reliability, and validity of the

480     arabic version of neck disability index in patients with neck pain. *Spine (Phila Pa 1976)*

481     2013;**38**:609–15. doi:10.1097/BRS.0b013e31828b2d09

482   28   Takeshita K, Hosono N, Kawaguchi Y, *et al.* Validity, reliability and responsiveness of the

483     Japanese version of the Neck Disability Index. *J Orthop Sci* 2013;**18**:14–21. doi:10.1007/s00776-

484     012-0304-y

485   29   Trouli MN, Vernon HT, Kakavelakis KN, *et al.* Translation of the Neck Disability Index and

486     validation of the Greek version in a sample of neck pain patients. *BMC Musculoskelet Disord*

487     2008;**9**:1–8. doi:10.1186/1471-2474-9-106

488   30   Tuttle N, Laakso L, Barrett R. Change in impairments in the first two treatments predicts outcome

489     in impairments, but not in activity limitations, in subacute neck pain: An observational study. *Aust*

490     *J Physiother* 2006;**52**:281–5. doi:10.1016/S0004-9514(06)70008-3

491   31   Ngo Trung, Stupar Maja, Coˆteˊ Pierre, Boyle Eleanor, Shearer Heather. A study of the test –

492     retest reliability of the self-perceived general recovery and self-perceived change in neck pain

493     questions in patients with recent whiplash-associated disorders. 2010;:957–62.

494     doi:10.1007/s00586-010-1289-x

495   32   Björklund M, Wiitavaara B, Heiden M. Responsiveness and minimal important change for the

496     ProFitMap-neck questionnaire and the Neck Disability Index in women with neck–shoulder pain.

497     *Qual Life Res* 2017;**26**:161–70. doi:10.1007/s11136-016-1373-8

498   33   Evans R, Bronfort G, Maiers M, *et al.* '" I know it " s changed '': a mixed-methods study of the

499     meaning of Global Perceived Effect in chronic neck pain patients. 2014;:888–97.

500     doi:10.1007/s00586-013-3149-y

501   34   Jorritsma W, Dijkstra PU, De Vries GE, *et al.* Detecting relevant changes and responsiveness of

502     Neck Pain and Disability Scale and Neck Disability Index. *Eur Spine J* 2012;**21**:2550–7.

21

503        doi:10.1007/s00586-012-2407-8

504    35    Monticone M, Frigau L, Vernon H, *et al.* Responsiveness and minimal important change of the

505        NeckPix© in subjects with chronic neck pain undergoing rehabilitation. *Eur Spine J*

506        2018;**27**:1324–31. doi:10.1007/s00586-017-5343-9

507    36    Monticone M, Ambrosini E, Vernon H, *et al.* Responsiveness and minimal important changes for

508        the Neck Disability Index and the Neck Pain Disability Scale in Italian subjects with chronic neck

509        pain. *Eur Spine J* 2015;**24**:2821–7. doi:10.1007/s00586-015-3785-5

510    37    Young BA, Walker MJ, Strunce JB, *et al.* Responsiveness of the Neck Disability Index in patients

511        with mechanical neck disorders. *Spine J* 2009;**9**:802–8. doi:10.1016/j.spinee.2009.06.002

512    38    Farooq MN, Mohseni-Bandpei MA, Gilani SA, *et al.* Urdu version of the neck disability index: A

513        reliability and validity study. *BMC Musculoskelet Disord* 2017;**18**:1–11. doi:10.1186/s12891-017-

514        1469-5

515    39    Williams GN, Gangel TJ, Arciero RA, *et al.* Comparison of the single assessment numeric

516        evaluation method and two shoulder rating scales. Outcomes measures after shoulder surgery. *Am

517        J Sports Med* 1999;**27**:214–21. doi:10.1177/03635465990270021701

518    40    Schmitt J, Abbott JH. Global Ratings of Change Do Not Accurately Reflect Functional Change

519        Over Time in Clinical Practice. *J Orthop Sport Phys Ther* 2015;**45**:106–11.

520        doi:10.2519/jospt.2015.5247

521    41    Chiarotto A, Ostelo RW, Boers M, *et al.* A systematic review highlights the need to investigate the

522        content validity of patient-reported outcome measures for physical functioning in patients with

523        low back pain. *J Clin Epidemiol* 2018;**95**:73–93. doi:10.1016/j.jclinepi.2017.11.005

524    42    Ailliet L, Knol DL, Rubinstein SM, *et al.* Definition of the construct to be measured is a

525        prerequisite for the assessment of validity. the Neck Disability Index as an example. *J Clin

526        Epidemiol* 2013;**66**:775-782.e2. doi:10.1016/j.jclinepi.2013.02.005

527

528

22

1
2
3      529
4
5      530
6
7      531
8
9      532
10
11
12     533      **Figure 1.** Flow diagram of included studies
13
14     534      **Figure 2**. Meta-analysis of Pearson's correlation coefficients between neck disability change scores and
15     535      GROC scores in patients with neck disorders based on 5 very good to excellent quality studies.
16
17     536      **Figure 3**. Meta-analysis of Spearman's correlation coefficients between neck disability change scores
18     537      and GROC scores in patients with neck disorders based on 6 very good to excellent quality studies.
19
20     538      **Figure 4**. Random effects univariate meta-regression between age and the Fisher's Z estimates. Each circle
21     539      represents a study and the size of the circle indicates the influence of that study on the model. The
22     540      regression prediction is illustrated by the straight line and the curved lines represent the 95% confidence
23     541      intervals. Age explained 68% of the variance in the model ($R^2$=0.68)
24
25
26     542
27
28     543
29
30     544
31     545
32
33     546
34
35     547
36
37     548
38
39     549
40
41     550
42
43     551
44
45     552
46
47     553
48     554
49
50     555
51
52     556
53
54     557
55
56     558
57
58                                                                                                              23
59
60

559

560

**Table 1**. Study Characteristics

| Study | Population | Setting | Sample Size | Properties Evaluated | GROC evaluated (ranked categories) | Interval |
|---|---|---|---|---|---|---|
| Bjorklund et al (2017) | Women with non-specific neck-shoulder pain | Not specified | 104 | Validity (correlation) Between NDI and GRoC | GRoC (7) 1. Very much worse; 2. Much worse; 3. Minimally worse; 4. No change; 5. Minimally improved; 6. Much improved; 7. Very much improved. | GRoC scale administered only after intervention at one time point (1 week) |
| Cleland et al (2006) | Patients with cervical radiculopathy | Hospital | 38 | Validity (correlation) Between NDI and GRoC Between PSFS and GRoC | GRoC (15) -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at follow up. Within a week over the period of 7 weeks. |
| Cleland et al. (2008) | Patients with neck pain only | 5 Outpatient physical therapy clinics | 137 | Validity (correlation) Between NDI and GRoC Between NPRS and GRoC | GRoC (15) -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at follow up. Within a week |
| Cook et al (2014) | Patients with any neck pain | Academic locations in Northeast Ohio | 56 | ROC curves and AUC to measure sensitivity and specificity. Binomial logistic regression analysis was also calculated to determine overall effect. | GRoC (15) -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | Baseline and at follow up 48- and 72 hours post baseline |
| Farooq et al. (2017) | Patients with neck pain | Physical therapy clinics | 106 | Validity (correlation) Between NDI-U and GRoC | GRoC (15) -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | GRoC was completed at three weeks after intervention |
| Guzy et al. (2013) | Patients with neck pain | Outpatient rehabilitation clinic | 95 | Validity (correlation) Between NDI-P and GRoC | GRoC (7) ''complete recovery'' over ''no change'' to ''my complaints are worse than ever'' | GRoC scale was completed at 2 weeks and at 4 weeks |
| Jorritsma et al. (2012) | Patients with chronic non-specific neck pain | Tertiary university center for rehabilitation | 76 | Validity (correlation) Between NDI and GRoC Between NPAD and GRoC | GPE (7) 3 (completely recovered) to zero (no change) to -3 (worse than ever) | After completion of the program varying from 3 to 5 months patients filled the GPE |

561

562

24

| | | | | | | |
|---|---|---|---|---|---|---|
| Kamper et al. (2010) | Patients with any whiplash-associated disorder. | Physical therapy clinics | 134 | Test-retest reliability | GPE (11) -5 (vastly worse) to zero (unchanged) to +5 (completely recovered) | Baseline, 6 weeks, and 12 months |
| Monticone et al. 2017 | Patients with chronic neck pain | Outpatient Rehabilitation Unit | 153 | Validity (correlation) Between NeckPix and GPE | GPE (5) (helped a lot = 1, helped = 2), one no change level (helped only a little = 3), and two worsening levels (did not help = 4, made things worse = 5) | At the end of treatment (8 weeks) and one year before follow-up |
| Monticone et al. 2015 | Patients with chronic neck pain | Outpatient Rehabilitation Unit | 200 | Validity (correlation) Between NDI and GPE Between NPDS and GPE | GPE (5) (helped a lot = 1, helped = 2), one no change level (helped only a little = 3), and two worsening levels (did not help = 4, made things worse = 5) | At the end of treatment 8 week |
| Ngo et al. (2010) | Patients with WAD. Most participants (69.6%) had grade II WAD. | Interviewed by person or by telephone in Ontario | 46 | Test-retest reliability | GPE (7) 1. General recovery question Completely better Much improved Slightly improved No change Slightly worse Much worse Worse than ever 2. Change in neck pain question: very much better, better, slightly better, no change, slightly worse, worse, or very much worse | 3-5 days |
| Shaheen et al. (2015) | Patients with neck pain lasting more than 3 months | 3 primary health centers | 70 | Validity (correlation) Between NDI-Ar and GRoC | GRoC (15) -7 (a very great deal worse) to zero (about the same) to +7 (a very great deal better) | 1 week |
| Takeshita et al. (2014) | Patients with neck pain, cervical radiculopathy and/or cervical myelopathy | Variety of clinics and hospital settings | 130 | Validity (correlation) Between NDI-J and GRoC | PGIC (7) much better, better, slightly better, unchanged, slightly worse, worse and much worse | Over 8 weeks |
| Trouli et al. (2008) | Patients with neck pain | Primary healthcare clinic | 68 | Validity (correlation) Between NDI-Gr and GRoC | GRoC (15) -7 (a very great deal worse) to -1 (almost the same, hardly any worse at all) and from 7 (a very great deal better) to 1 (almost the same, hardly any better at all) | Within 2 months but 1 week for test-retest |

25

| | | | | | | |
|---|---|---|---|---|---|---|
| Tuttle et al. (2006) | Patients with neck pain for more than 2 weeks | Private physiotherapy clinics | 29 | Validity (correlation) Between NDI and GPE Between PSFS and GPE Between VAS and GPE Between ROM and GPE | GPE (11) −5 is vastly worse and +5 is completely recovered | 6 weeks |
| Young et al. (2009) | Patients presenting with mechanical neck pain | Outpatient physical therapy clinics. | 91 | Validity (correlation) | GRoC (15) -7 (''a very great deal worse'') to 0 (''about the same'') to +7 (''a very great deal better'') | 3 weeks |

563  NDI = Neck Disaiblity Index, NPRS=Numeric Pain Rating Scale, PSFS= Patient Specific Functional Scale, ROC= Receiver Operator
564  Characteristic, VAS=Visual Analog Scale, NPAD=Neck Pain and Disability Scale, AUC= Area Under the Curve, ROM=Range of
565  Motion

26

567 **TABLE 2.** Summary of Psychometric Properties Reported in Studies and COSMIN Risk of Bias (RoB)
568 and Quality studies

| Study | Psychometric Properties Reported | COSMIN RoB | COSMIN Rating*§ (Criteria) | Quality of Studies** (QACMRR) |
|---|---|---|---|---|
| Bjorklund et al (2017) | Validity (correlation) | Very Good | ? | Excellent |
| Cleland et al (2006) | Validity (correlation) | Very Good | + | Excellent |
| Cleland et al. (2008) | Validity (correlation) | Very Good | - | Excellent |
| Cook et al (2014) | Sensitivity Specificity | Very Good Very Good | + | Excellent |
| Farooq et al. (2017) | Validity (correlation) | Very Good | + | Excellent |
| Guzy et al. (2013) | Validity (correlation) | Very Good | ? | Very good |
| Jorritsma et al. (2012) | Validity (correlation) | Very Good | ? | Excellent |
| Kamper et al. (2010) | Test-retest reliability | Very Good | + | Excellent |
| Monticone et al. (2017) | Validity (correlation) | Very Good | ? | Excellent |
| Monticone et al. (2015) | Validity (correlation | Very Good | ? | Excellent |
| Ngo et al. (2010) | Test-retest reliability | Very Good | + | Excellent |
| Shaheen et al. (2015) | Validity (correlation) | Very Good | ? | Excellent |
| Takeshita et al. (2014) | Validity (correlation) | Very Good | ? | Very good |
| Trouli et al. (2008) | Validity (correlation) | Very Good | + | Excellent |
| Tuttle et al. (2006) | Validity (correlation) | Very Good | ? | Excellent |
| Young et al. (2009) | Validity (correlation) | Very Good | ? | Excellent |

569 COSMIN, Consensus-based Standards for the Selection of health Measurement Instruments, Criteria for good measurement
570 properties: '+' sufficient; '-'insufficient; '?' indeterminate. §§ The grading for the quality of the evidence based on the modified
571 GRADE approach is not applicable. **Quality Appraisal for Clinical Measurement Research Reports Evaluation Form
572 (QACMRR).

573

574

575

576

577

578

579

580

581

27

582 **TABLE 3**. Quality Appraisal for Clinical Measurement Research Reports Evaluation Form

| Study | Item Evaluation Criteria* | | | | | | | | | | | | Total (%) | Quality Summary |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | **1** | **2** | **3** | **4** | **5** | **6** | **7** | **8** | **9** | **10** | **11** | **12** | | |
| Bjorklund et al (2017) | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Cleland et al. (2008) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Trouli et al. (2008) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Tuttle et al. (2006) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Kamper et al. (2010) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 96 | Excellent |
| Cook et al (2014) | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 92 | Excellent |
| Jorritsma et al. (2012) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Cleland et al (2006) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Monticone et al. (2017) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Monticone et al. (2015) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Ngo et al. (2010) | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 1 | 2 | 92 | Excellent |
| Shaheen et al. (2013) | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 92 | Excellent |
| Farooq et al. (2017) | 2 | 2 | 1 | 2 | 2 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 92 | Excellent |
| Young et al. (2009) | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 92 | Excellent |
| Guzy et al. (2013) | 2 | 2 | 1 | 2 | 1 | 2 | 2 | 2 | 1 | 2 | 2 | 2 | 88 | Very good |
| Takeshita et al. (2014) | 2 | 2 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 88 | Very good |

583 *Item Evaluation Criteria: 1. Thorough literature review to define the research question; 2. Specific inclusion/exclusion

584 criteria; 3. Specific hypotheses; 4. Appropriate scope of psychometric properties; 5. Sample size; 6. Follow-up; 7. The

585 authors referenced specific procedures for administration, scoring, and interpretation of procedures; 8. Measurement

586 techniques were standardized; 9. Data were presented for each hypothesis; 10. Appropriate statistics-point estimates; 11.

587 Appropriate statistical error estimates; 12. Valid conclusions and clinical recommendations.

28

588    *Total score = (sum of subtotals ÷ 24 × 100). If for a specific paper an item is deemed NA (Not Applicable), then, Total score*

589    *= (sum of subtotals ÷ (2 × number of Applicable items) × 100).*

590    *NA – Not Applicable. The subsections no. 6, asks for percentage of retention/follow up. This subsection only applies to*

591    *reliability test-retest studies*

592    *Quality Summary: Poor (0%-30%), Fair (31%-50%), Good (51%-70%), Very good (71%-90%), Excellent (>90%):*

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

29

617 **TABLE 4**. Summary of reliability properties of GRoC scales

| Study | Type of Reliability | Reliability Estimates | COSMIN | Quality of Studies |
|-------|---------------------|----------------------|--------|--------------------|
| Kamper et al. (2010) | Test-retest | Intra-class correlation coefficients (ICC) <br> 0.99 (0.99 – 0.99) – baseline <br> 0.96 (0.95 – 0.97) – at six weeks <br> 0.92 (0.89 – 0.94) at twelve months. | Very Good | Excellent |
| Ngo et al. (2010) | Test-retest | Intra-class correlation coefficients (ICC) <br> 0.70 (0.60–0.80) – at six weeks (General recovery) <br> 0.80 (0.72–0.87) – at six weeks (neck pain questions) <br><br> Weighted Kappa <br> 0.70 (0.42–0.98) – at six weeks (General recovery) <br> 0.80 (0.51–1.0) – at six weeks (neck pain questions) <br><br> Dichotomized response options for recovery (K statistics) <br> 0.85 (0.64–1) when ''recovered'' was defined ''completely better' <br> 0.81 (0.64–0.99) when defined as ''completely better'' or ''much improved <br><br> Dichotomized response options for change in neck pain questions (K statistics) <br> 0.46 (0.20–0.74) when ''recovered'' was defined as ''very much better'' <br> 0.80 (0.62–0.99) when defined as ''very much better'' or ''better' <br><br> Recall questions (K statistics) <br> the kappa coefficient was 1 for participants who remembered their previous answers to the general recovery question; 0.88 (0.64–1) for those who did not remember and 0.50 (0.02– 0.98) for participants who were not asked the question. <br><br> The kappa coefficient was 1 for participants who remembered their previous answers to the change in neck pain question; 0.74 (0.41–1) for those who did not remember and 0.66 (0.22–1) for participants who were not asked the question. | Very Good | Excellent |

618

619

620

621

622

623

624

625

626

30

627 **TABLE 5**. Summary of validity properties of GRoC scales

| Study | Type of Reliability | Validity Estimates | COSMIN | Quality of Studies |
|---|---|---|---|---|
| Bjorklund et al (2017) | Spearman's correlation between the change scores of GRoC and ProFitMap-neck<br><br>GRoC and NDI | rho = 0.47, (p<0.05)<br>rho = 0.59, (p<0.05) | Very Good | Excellent |
| Cleland et al. (2006) | Correlations (Pearson r) between change scores NDI and GRoC<br>PSFS and GRoC | r = 0.19<br><br>r = 0.82 | Very Good | Excellent |
| Cleland et al. (2008) | Correlations (Pearson r) between change scores NDI and GRoC<br>NRS and GRoC | r = 0.58<br>r = 0.57 | Very Good | Excellent |
| Cook et al. (2014) | Receiver operator characteristics (ROC) Within-session change Between-session change<br><br>Between session change of Pain and GROC Sensitivity Specificity | AUC = 0.61<br>AUC = 0.76, >36.7% change in pain<br><br>Odds ratio = 7.3 (2.1, 24.7)<br>65.6% (57.9, 74.6)<br>79.2% (62.2, 91.1) | Very Good | Excellent |
| Farooq et al. (2017) | Correlations (Pearson r) NDI-U | r =0.50 | Very Good | Excellent |
| Guzy et al. (2013) | Correlations (Pearson r) NDI vs GROC | Two- week interval (r = -0.73)<br>Four-week interval (r = -0.56) | Very Good | Very good |
| Jorritsma et al. (2012) | Correlation between change scores of NPAD and GPE | r = 0.49 (95 % CI 0.30–0.64) | Very Good | Excellent |
| Monticone et al. (2017) | Correlations (Spearman) between change scores of the NeckPix© and GPE | rho = 0.69–0.82 | Very Good | Excellent |
| Monticone et al. (2015) | Correlation (Spearman) between change scores NDI-I and GPE<br>NDPS and GPE | rho = 0.71, p<0.01<br>rho = 0.59, p<0.01 | Very Good | Excellent |
| Shaheen et al. (2013) | Correlations (Spearman's) NDI-Ar and GROC | rho = 0.81, p<o.oo1 | Very Good | Excellent |
| Takeshita et al. (2014) | Correlations NDI and PGIC<br>NDI-J and PGIC | Spearman (rho)<br>rho = 0.47, p<o.oo1<br>rho = 0.59, p<o.oo1 | Very Good | Very good |
| Trouli et al. (2008) | Correlation (Spearman's) GROC vs Gr-NDI | rho = 0.30, p=0.02 | Very Good | Excellent |
| Tuttle et al. (2006) | Correlations (Spearman's) NDI vs GPE (post 1, minus pre-1) NDI vs GPE (post 2, minus pre-1) NDI vs GPE (post 2, minus pre-2)<br><br>PSFS vs GPE (post 1, minus pre-1) PSFS vs GPE (post 2, minus pre-1) PSFS vs GPE (post 2, minus pre-2) | rho = 0.17<br>rho = 0.01<br>rho = 0.03<br><br>rho = 0.06<br>rho = 0.03<br>rho = 0.03 | Very Good | Excellent |

31

| | | | | |
|---|---|---|---|---|
| | Pain Intensity (post 1, minus pre-1) | rho = 0.00 | | |
| | Pain Intensity (post 2, minus pre-1) | rho = 0.05 | | |
| | Pain Intensity (post 2, minus pre-2) | rho = 0.01 | | |
| | Total ROM (post 1, minus pre-1) | rho = 0.03 | | |
| | Total ROM (post 2, minus pre-1) | rho = 0.01 | | |
| | Total ROM (post 2, minus pre-2) | rho = 0.00 | | |
| Young et al. (2009) | Correlations (Pearson's) between change scores NDI and GRoC | r =0.52 (p<0.01) | Very Good | Excellent |
| Monticone et al. (2015) | Correlation (Spearman) between change scores NDI-I and GPE NDPS and GPE | rho = 0.71, p<0.01 rho = 0.59, p<0.01 | Very Good | Excellent |

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

32

650

651    **Box 1.** Questions of Global Rating of Change (GROC) scales

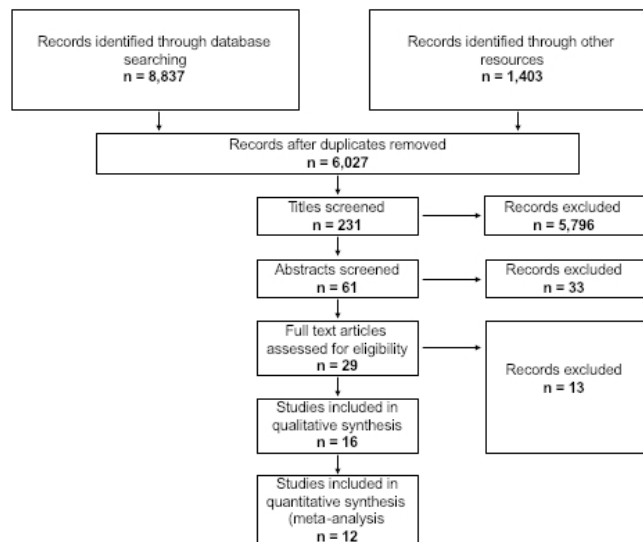| Author | GROC (ranked categories) | Patients with neck disorders were asked: |
|---|---|---|
| Bjorklund et al. (2017) | GROC (7) | *"Compared to before the treatment of the study started, my overall status is now"* <br><br> *"Compared to before the treatment of the study started, my status regarding my neck–shoulder problem is now"* |
| Evans et al (2014) | GPE (9) | *"Overall, how much has your neck pain changed since you started treatment in the study?"* |
| Kamper et al. (2010) | GPE (11) | *"With respect to your whiplash injury how would you describe yourself now compared to immediately after your accident"* |
| Monticone et al. (2017) | GPE (5) | *"Overall, how much did the treatment you received help your fear of movement due to current neck pain?* <br><br> *"Overall, how much did the treatment you delivered help your subject's fear of movement due to her/ his current neck pain?"* |
| Monticone et al. (2015) | GPE (5) | *"Overall, how much did the treatment you received help your neck problem?"* |
| Ngo et al. (2010) | GPE (7) | *"How well do you feel you are recovering from your injuries?"* <br><br> *"How do you feel your neck pain has changed since the injury?"* |

652

653

654

655

656

33

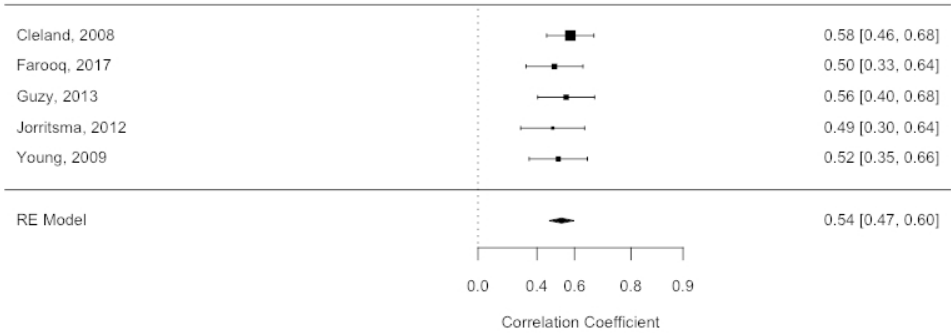Figure 1. Flow diagram of included studies

60x34mm (300 x 300 DPI)

Figure 2. Meta-analysis of Pearson's correlation coefficients between neck disability change scores and GROC scores in patients with neck disorders based on 5 very good to excellent quality studies.
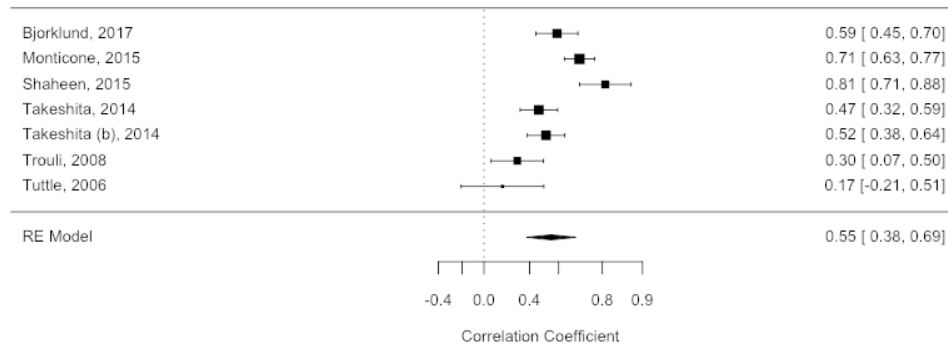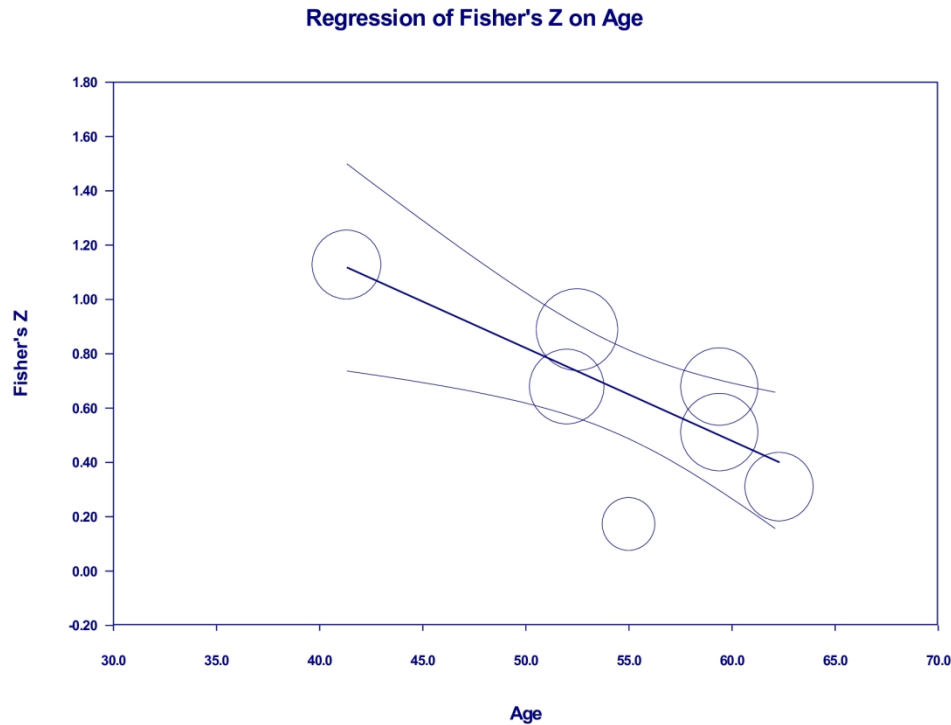
67x34mm (300 x 300 DPI)

Figure 3. Meta-analysis of Spearman's correlation coefficients between neck disability change scores and GROC scores in patients with neck disorders based on 6 very good to excellent quality studies.

67x34mm (300 x 300 DPI)

Figure 4. Random effects univariate meta-regression between age and the Fisher's Z estimates. Each circle represents a study and the size of the circle indicates the influence of that study on the model. The regression prediction is illustrated by the straight line and the curved lines represent the 95% confidence intervals. Age explained 68% of the variance in the model (R2=0.68)

160x118mm (300 x 300 DPI)

**Appendix 1**

**Search terms**

MEDLINE-OVID
1. exp "outcome and process assessment (health care)"/ or "outcome assessment (health care)"/
or treatment outcome/
2. outcome?.ti.
3. exp "Range of Motion, Articular"/
4. Pain Measurement/
5. exp disability evaluation/
6. "Recovery of Function"/
7. Questionnaires/
8. self-report.tw.
9. ((impairment or disability or function) adj2 (measure? or scale? or evaluation?)).tw.
10. range of motion.tw.
11. (strength adj2 (measure? or scale? or evaluation?)).tw.
12. (outcome? adj2 (measure* or scale? or indicator?)).tw.
13. or/1-12
14. "reproducibility of results"/
15. exp "Sensitivity and Specificity"/
16. reliability.mp.
17. validity.mp.
18. responsiveness.mp.
19. Psychometrics/
20. rasch.mp.
21. factor analysis, statistical/
22. factor analysis.tw.
23. differential functioning.mp.
24. (validity or validation).mp. [mp=title, original title, abstract, name of substance word, subject
heading word, unique identifier]
25. (validity or validation).mp.
26. item difficulty.mp.
27. translation.tw.
28. or/14-27
29. 13 and 28
30. Neck Pain/
31. exp Brachial Plexus Neuropathies/
32. exp neck injuries/ or exp whiplash injuries/
33. cervical pain.mp.
34. neckache.mp.
35. whiplash.mp.
36. cervicodynia.mp.
37. cervicalgia.mp.
38. brachialgia.mp.
39. brachial neuritis.mp.

40. brachial neuralgia.mp.

41. neck pain.mp.

42. neck injur*.mp.

43. brachial plexus neuropath*.mp.

44. brachial plexus neuritis.mp.

45. thoracic outlet syndrome/ or cervical rib syndrome/

46. Torticollis/

47. exp brachial plexus neuropathies/ or exp brachial plexus neuritis/

48. cervico brachial neuralgia.ti,ab.

49. cervicobrachial neuralgia.ti,ab.

50. (monoradicul* or monoradicl*).tw.

51. or/30-50

52. exp headache/ and cervic*.tw.

53. exp genital diseases, female/

54. genital disease*.mp.

55. or/53-54

56. 52 not 55

57. 51 or 56

58. neck/

59. neck muscles/

60. exp cervical plexus/

61. exp cervical vertebrae/

62. atlanto-axial joint/

63. atlanto-occipital joint/

64. Cervical Atlas/

65. spinal nerve roots/

66. exp brachial plexus/

67. (odontoid* or cervical or occip* or atlant*).tw.

68. axis/ or odontoid process/

69. Thoracic Vertebrae/

70. cervical vertebrae.mp.

71. cervical plexus.mp.

72. cervical spine.mp.

73. (neck adj3 muscles).mp.

74. (brachial adj3 plexus).mp.

75. (thoracic adj3 vertebrae).mp.

76. neck.mp.

77. (thoracic adj3 spine).mp.

78. (thoracic adj3 outlet).mp.

79. trapezius.mp.

80. cervical.mp.

81. cervico*.mp.

82. 80 or 81

83. exp genital diseases, female/

84. genital disease*.mp.

85. exp *Uterus/

86. 83 or 84 or 85
87. 82 not 86
88. 58 or 59 or 60 or 61 or 62 or 63 or 64 or 65 or 66 or 67 or 68 or 69 or 70 or 71 or 72 or 73 or 74 or 75 or 76 or 77 or 78 or 79 or 87
89. exp pain/
90. exp injuries/
91. pain.mp.
92. ache.mp.
93. sore.mp.
94. stiff.mp.
95. discomfort.mp.
96. injur*.mp.
97. neuropath*.mp.
98. or/89-97
99. 88 and 98
100. Radiculopathy/
101. exp temporomandibular joint disorders/ or exp temporomandibular joint dysfunction syndrome/
102. myofascial pain syndromes/
103. exp "Sprains and Strains"/
104. exp Spinal Osteophytosis/
105. exp Neuritis/
106. Polyradiculopathy/
107. exp Arthritis/
108. Fibromyalgia/
109. spondylitis/ or discitis/
110. spondylosis/ or spondylolysis/ or spondylolisthesis/
111. radiculopathy.mp.
112. radiculitis.mp.
113. temporomandibular.mp.
114. myofascial pain syndrome*.mp.
115. thoracic outlet syndrome*.mp.
116. spinal osteophytosis.mp.
117. neuritis.mp.
118. spondylosis.mp.
119. spondylitis.mp.
120. spondylolisthesis.mp.
121. or/100-120
122. 88 and 121
123. exp neck/
124. exp cervical vertebrae/
125. Thoracic Vertebrae/
126. neck.mp.
127. (thoracic adj3 vertebrae).mp.
128. cervical.mp.
129. cervico*.mp.

130. 128 or 129
131. exp genital diseases, female/
132. genital disease*.mp.
133. exp *Uterus/
134. or/131-133
135. 130 not 134
136. (thoracic adj3 spine).mp.
137. cervical spine.mp.
138. 123 or 124 or 125 or 126 or 127 or 135 or 136 or 137
139. Intervertebral Disk/
140. (disc or discs).mp.
141. (disk or disks).mp.
142. 139 or 140 or 141
143. 138 and 142
144. herniat*.mp.
145. slipped.mp.
146. prolapse*.mp.
147. displace*.mp.
148. degenerat*.mp.
149. (bulge or bulged or bulging).mp.
150. 144 or 145 or 146 or 147 or 148 or 149
151. 143 and 150
152. intervertebral disk degeneration/ or intervertebral disk displacement/
153. intervertebral disk displacement.mp.
154. intervertebral disc displacement.mp.
155. intervertebral disk degeneration.mp.
156. intervertebral disc degeneration.mp.
157. 152 or 153 or 154 or 155 or 156
158. 138 and 157
159. 57 or 99 or 122 or 151 or 158
160. animals/ not (animals/ and humans/)
161. 159 not 160
162. exp *neoplasms/
163. exp *wounds, penetrating/
164. 162 or 163
165. 161 not 164
166. 29 and 165
167. guidelines as topic/
168. practice guidelines as topic/
169. guideline.pt.
170. practice guideline.pt.
171. (guideline? or guidance or recommendations).ti.
172. consensus.ti.
173. or/167-172
174. meta-analysis/
175. exp meta-analysis as topic/

176. (meta analy* or metaanaly* or met analy* or metanaly*).tw.
177. review literature as topic/
178. (collaborative research or collaborative review* or collaborative overview*).tw.
179. (integrative research or integrative review* or intergrative overview*).tw.
180. (quantitative adj3 (research or review* or overview*)).tw.
181. (research integration or research overview*).tw.
182. (systematic* adj3 (review* or overview*)).tw.
183. (methodologic* adj3 (review* or overview*)).tw.
184. exp technology assessment biomedical/
185. (hta or thas or technology assessment*).tw.
186. ((hand adj2 search*) or (manual* adj search*)).tw.
187. ((electronic adj database*) or (bibliographic* adj database*)).tw.
188. ((data adj2 abstract*) or (data adj2 extract*)).tw.
189. (analys* adj3 (pool or pooled or pooling)).tw.
190. mantel haenszel.tw.
191. (cohrane or pubmed or pub med or medline or embase or psycinfo or psyclit or psychinfo or psychlit or cinahl or science citation indes).ab.
192. or/174-191
193. 173 or 192
194. 166 and 193

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

**Quality Appraisal for Clinical Measurement Research Reports**

**Evaluation Form**

Authors: _____ Year: _____ Rater: _____

*Use this form to rate the quality of a clinical measurement study. To decide which score to provide for each item on your quality checklist, pick the descriptor that sounds <u>most</u> like what was reported in the study you are evaluating.  Items rank descriptors are provided in the guide. (Forms and guides to extract study data for evidence synthesis are available from developer at macderj@mcmaster.ca)*

| Evaluation criteria | Score | | |
|---|---|---|---|
| **Study question** | 2 | 1 | 0 |
| 1. Was the relevant background work cited to define what is currently known about the measurement properties of measures under study, and the potential contributions of the current research question to informing that knowledge base? | | | |
| **Study Design** | | | |
| 2. Were appropriate inclusion/exclusion criteria defined? | | | |
| 3. Were specific clinical measurement questions/hypotheses identified? | | | |
| 4. Was an appropriate scope of measurement properties considered? | | | |
| 5. Was an appropriate sample size used? | | | |
| 6. Was appropriate retention/follow-up obtained? (for studies involving retesting; otherwise  n/a) | | | |
| **Measurements** | | | |
| 7. Were specific descriptions provided of the measure under study and the method(s) used to administer it? | | | |
| 8.  Were standardized procedures used to administer all study measures in a manner that minimized potential sources of error/bias (including the study measure and its comparators)? | | | |
| **Analyses** | | | |

| | | | |
|---|---|---|---|
| 9. Were analyses conducted for each specific hypothesis or purpose? | | | |
| 10. Were appropriate statistical tests performed to obtain point estimates of the measurement properties? | | | |
| 11. Were appropriate ancillary analyses done to quantify the confidence in the estimates of the clinical measurement property (Precision/Confidence intervals; benchmark comparisons/ROC curves, alternate forms of analysis like SEM/MID, etc.)? | | | |
| **Recommendations** | | | |
| 12. Were clear, specific and accurate conclusions made about the clinical measurement properties; that were associated with appropriate clinical measurement recommendations and supported by the study objectives, analysis and results? | | | |
| **Subtotals** (of columns 1 and 2) | | | |
| **Total score** (sum of subtotals/24*100); if for a specific paper or topic an item is deemed inappropriate then you can sum of items/2*number of items *100 | | | |

© MacDermid 2011

### Quality Appraisal of a Clinical Measurement Study

### Interpretation Guide

To decide which score to provide for each item on your quality checklist, read the following descriptors. Pick the descriptor that sounds _most_ like the study you were evaluating with respect to a given item. If there is no documentation about any specific aspect of an item; then you must evaluate assuming that it was not done. Given the diversity in clinical measurement properties and design options, the evaluator has to make judgments using the criteria below and extend the principles to specific aspects that may not be covered in these brief exemplars. In many cases, the study will not look exactly like the descriptor so there will be some interpretation as to which level of optimal methods for clinical measurement studies have been achieved. In such cases, the evaluator can use the general approach that if this study research design and conduct is consistent with best practice (score=2); is acceptable but suboptimal (score=1); is not done/documented, substantially inadequate or inappropriate (score=0).

| | **Descriptors** | |
|---|---|---|
| **Study question** | | |
| Score | | |
| 1 | 2 | The authors:<br><br>- performed a thorough literature review indicating what is currently known, and not known, about the clinical measurement properties of the instruments or tests under study<br>- presented a critical, and unbiased view of what is known about the current measurement properties<br>- indicated how the current research question fills a gap in the current knowledge base<br>- established a research question based on the above. |
| | 1 | All of the above criteria were not fulfilled, but a sound rationale was provided for the research question. |
| | 0 | A foundation for the current research question was not clear; and the rationale was not founded on previous literature. |
| **Study design** | | |

| 2 | 2 | Specific inclusion/exclusion criteria for the study were defined, that described the patients enrolled. The subjects were described in terms of health condition/demographics, key relevant outcome mediators and the recruitment context (setting). |
|---|---|---|
| | 1 | Some information on participants and place is provided (not all of above). For example, age/sex/diagnosis and the name or type of the practice is listed; but no additional information. |
| | 0 | No information on type of clinical settings or study participants is provided (other than number/mean age). |
| 3 | 2 | Specific hypotheses or research questions are provided. The stated study purpose provides specific research questions or hypotheses that indicate which specific measurement properties will be evaluated. This should include the specific type of reliability (intra/inter-rater or test-retest) being tested or the type of validity (construct/criterion/content; longitudinal/concurrent; convergent/divergent) being tested. A prior hypothesis should describe the level of reliability expected; and for validity, expected relationships (strength of associations) or constructs. |
| | 1 | The types of reliability and validity being tested were apparent in the methods/title, but clear and specific research questions or hypotheses were not specified. |
| | 0 | Specific types of reliability or validity under evaluation were not clearly defined nor were specific hypotheses on reliability and validity stated. ("*The purpose of this study was to investigate the reliability and validity of*…" can be rated as zero if no further detail on the types of reliability and validity or the nature of specific hypotheses is stated). |
| 4 | 2 | An appropriate scope of clinical measurement properties would be indicated by <br><br> 1. A detailed focus on reliability that included multiple forms of reliability (at least two of – intra-rater, inter-rater, test retest); as well as both relative and absolute reliability (e.g., ICCs and SEM/MID or limits of agreement) <br> 2. A detailed focus on validity that included multiple forms of validity (content (judgmental); structured (e.g., expert review/survey, qualitative interviews, ICF linking) or structural (e.g., factor analyses or Rasch), construct (known group differences; convergent/divergent associations), criterion (concurrent/predictive), responsiveness; predictive, evaluative or discriminative properties were established <br> 3. Three or more indicators of reliability and validity were examined concurrently and provide a rich view on measurement properties. |
| | 1 | Two or more clinical measurement properties were evaluated, however, scope was narrow and did not meet above criteria. (e.g., internal consistency and one other indicator of validity or reliability ). |
| | 0 | The scope of clinical measurement properties was very narrow as indicated by a narrow evaluation of only one form of reliability or validity. |

| 5 | 2 | Authors performed a sample size calculation and obtained their recruitment targets. Post-doc power analyses and/or confidence intervals confirm that the sample size was sufficient to define relatively precise estimates of reliability or validity. |
|---|---|---|
|   | 1 | The authors provide an acceptable rationale for the number of subjects included in the study, but did not present specific sample size calculations or post-doc power analyses (or had a sample >100 but no justification). |
|   | 0 | Size of the sample was not rationalized or is clearly underpowered. |
| 6 | 2 | 90% or more of the patients enrolled for study were re-evaluated. |
|   | 1 | 70% or more of the enrolled patients were re-evaluated. |
|   | 0 | Less than 70% of the patients enrolled in the study were re-evaluated |
| **Measurements** | | |
| 7 | 2 | Documentation is provided for how the studied test is performed.  This includes adequate description of the measure/test and how it is administered or scored. The authors may provide or reference a published manual/article that outlines specific procedures for administration, scoring (including scoring algorithms, handling of missing data) and interpretation that included any necessary information about positioning/active participation of the client, any special equipment required, calibration of equipment if necessary, training required, cost, examiner procedures/actions. If no manual is available, then the text describes key details of procedures in sufficient detail so they could be replicated. |
|   | 1 | The test(s) and its administration procedures are referenced; but there is inadequate description of the test procedures. |
|   | 0 | Minimal description of test procedures without appropriate references. |

| 8 | 2 | This item addresses the overall study procedures for administering all study measures (study measure and its comparators) in an unbiased way. Test procedures should not introduce systematic errors in the estimation of the clinical measurement properties. This includes standardized procedures for who completed or administered the measures. For self-report, this includes order of presentation, who completed at what time interval; handling of missing items. If relevant, then the paper should include how cultural literacy issues were handled (e.g., exclusion, assisted or surrogate completion). For impairment measures, procedures would include calibration of any equipment; use of consistent measurement tools and scoring, a priori exclusion of any participants likely to give invalid results/unable to complete testing (not exclusion of after enrollment); use of standardized instructions and test procedures.  This can include order of administration of test and quality checking of scores.  For reliability testing, the appropriate retest interval will depend on the nature of the condition; but for acute conditions it may require retesting within 48 hours; whereas chronic/stable conditions are commonly retested within 4-14 days.  For estimation of clinical change, retest intervals should be ones during which a meaningful clinical change would have occurred (and from an intervention with known effectiveness). The evaluator decides overall whether this has sufficiently been addressed by the methods described. |
|---|---|---|
| | 1 | No obvious sources of bias in the study test protocol or how tests were performed/administered is apparent; but there were suboptimal procedures or an inadequate description of the measurement protocol to be insured control of bias or that procedures were standardized. |
| | 0 | No description of the overall procedures for administering study tests; OR an obvious source of bias in data collection methods. |

**Analyses**

| 9 | 2 | Authors clearly defined which specific analyses were conducted for each of the stated specific hypotheses/questions of the study. This may be accomplished through organization of the results under specific subheadings or by demarcating which analyses addressed specific clinical measurement properties.  Data was presented for each hypothesis/research question posed. |
|---|---|---|
| | 1 | Data was presented that addressed each of the measurement questions posed, but authors did not link specific analyses to specific research questions or hypotheses. |
| | 0 | Data was not presented for every hypothesis or clinical measurement property outlined in the purposes or methods. |

| 10 | 2 | <u>Tests selected</u> - Appropriate statistical tests were conducted to calculate a point estimate for clinical measurement properties.  Examples are provided below; but are not exhaustive.

1.  Reliability (Relative=ICCs (Shrout & Fleiss, 1979) for quantitative, Kappa (Landis & Koch, 1977) for nominal data); absolute (SEM or plot of score differences vs. average score showing mean and  2SD limit – as per Altman and Bland) (Bland & Altman, 1986; Bland & Altman, 1987)

2.  Clinical relevance - minimal detectable change, clinically important difference (Jaeschke, Singer, & Guyatt, 1989; Beaton et al., 2001; Wells et al., 2001)

3.  Validity

a. Validity associations - Pearson correlations for normally distributed data, Spearman rank correlations for ordinal data; or other correlations, if appropriate

b. Validity tests of significant difference - an appropriate global test like analysis of variance was used where indicated, with post-hoc tests that adjusted for multiple testing

c. Validity of items scaling/responses - Rasch analysis or item response (Baylor et al., 2011; Pallant & Tennant, 2007; Kyngdon, 2006; Cipriani, Fox, Khuder, & Boudreau, 2005; Smith, Jr., Conrad, Chang, & Piazza, 2002)

4. Responsiveness (Beaton, Bombardier, Katz, & Wright, 2001)- standardized response means or effect sizes or other recognized responsiveness indices were used. |
| | 1 | Appropriate statistical tests were used in some instances; but suboptimal choices were made in other analyses. |
| | 0 | Inappropriate use of statistical tests - incorrect tests for type of data; or a lack of analysis |
| 11 | 2 | The study goes beyond a single statistical point estimate of a clinical measurement property and providing supporting statistical analyses that increases confidence in the findings in terms of precision of the (key) indicator; or provide an alternate form of analysis of the clinical measurement property. The evaluator decides if these analyses are appropriate and informative.  For example, with reliability, at least 2 of the following would constitute appropriate and informative analysis beyond a point estimate a reliability coefficient: 1. confidence intervals around the point estimate; 2. Comparison to appropriate referenced benchmarks or standards; or 3. SEM or MDC.  For correlations, tests of significance or confidence intervals were presented and indicators of the criterion benchmarks were provided.  For studies involving cross-cultural validation, the analyses should compare multiple clinical measurement properties previously established for the measure and explain the extent to which the translated version is in accordance with these previously reported properties on the source measure. |

| | 1 | Either precision definition (confidence intervals) or appropriate benchmark comparison were used - NOT both. OR Some analyses were associated with indicators of precision or alternate form of analysis -but not all key indicators. |
|---|---|---|
| | 0 | Inappropriate use of benchmarks or confidence intervals; or indicators of precision or alternate form are absent |

**Recommendations**

| 12 | 2 | Authors made specific conclusions and clinical measurement recommendations that were clearly related to each hypotheses/question posed in the study and that were supported by the data presented. Ideal recommendations would state the estimated status of the clinical measurement property, the confidence in the estimate and the context for which those apply. To achieve a 2, the conclusion must be specific; and conclusions cannot overstate the clinical measurement properties observed the study; nor ignore suboptimal measurement properties found. |
|---|---|---|
| | 1 | Authors made conclusions and clinical measurement recommendations that were basically true (supported by study data); but vague. That is, they do not specify the extent, confidence or context of the findings. (The measure is "reliable and valid ") OR authors made specific clinical measurement recommendations; but for only some of the study hypotheses. |
| | 0 | Authors did not make conclusions about clinical measurement; OR made recommendations that were in contradiction to the actual data presented |

© MacDermid 2011

**List with excluded studies with reasons**

| | |
|---|---|
| 1. Abbott et al 2014 | Ineligible population |
| 2. Beattie et al 2011 | Ineligible population (less than 50%) |
| 3. Hoeskstra et al 2014 | No properties for GRoC scales |
| 4. Chansirinukor 2019 | No properties for GRoC scales |
| 5. Chien et al 2015 | No properties for GRoC scales |
| 6. Cruz et al. 2015 | No properties for GRoC scales |
| 7. Foroutani et al 2018 | No English (Persian language) |
| 8. Gagnon et al 2018 | Ineligible population |
| 9. Hefford et al 2012 | Ineligible population |
| 10. Hung et al 2019 | Ineligible population |
| 11. Sharma et al 2017 | Ineligible population |
| 12. Stevens et al 2019 | Ineligible population |
| 13. Meyer et al 2014 | Ineligible population |

# 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 |
| **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. | 3-5 |
| Objectives | 4 | Provide an explicit statement of questions being addressed with reference to participants, interventions, comparisons, outcomes, and study design (PICOS). | 4-5 |
| **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. | 5 |
| 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 |
| 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. | 6 |
| Search | 8 | Present full electronic search strategy for at least one database, including any limits used, such that it could be repeated. | Appendix1 |
| 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). | 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. | 6-7 |
| Summary measures | 13 | State the principal summary measures (e.g., risk ratio, difference in means). | 8-9 |
| 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. | 8-9 |

# 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). | 8-9 |
| Additional analyses | 16 | Describe methods of additional analyses (e.g., sensitivity or subgroup analyses, meta-regression), if done, indicating which were pre-specified. | 8=9 |
| **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. | 9 |
| 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. | 9-10 |
| 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 |
| 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. | 10-12 |
| Synthesis of results | 21 | Present results of each meta-analysis done, including confidence intervals and measures of consistency. | 13 |
| Risk of bias across studies | 22 | Present results of any assessment of risk of bias across studies (see Item 15). | 10 |
| Additional analysis | 23 | Give results of additional analyses, if done (e.g., sensitivity or subgroup analyses, meta-regression [see Item 16]). | 13 |
| **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). | 14-15 |
| 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 |
| Conclusions | 26 | Provide a general interpretation of the results in the context of other evidence, and implications for future research. | 14-15 |
| **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. | 18 |

*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