2-12-2021

# Strengths and challenges of the COSMIN tools in the appraisal of outcome measures: A case example for speech-language therapy

Elaine Kwok
*McMaster University*, kwoke3@mcmaster.ca

Peter Rosenbaum
*McMaster University*

Nancy Thomas-Stonell
*McMaster University*

Barbara Jane Cunningham
*Western University*, bjcunningham@uwo.ca

1    **Title:** Strengths and challenges of the COSMIN tools in the appraisal of outcome measures: A

2                    case example for speech-language therapy

3    **Abstract:**

4    *Background:* The Consensus-based Standards for the selection of health Measurement

5    Instruments (COSMIN) is an international initiative that offers standardized and validated tools

6    to guide the appraisal of outcome measures in healthcare.

7    *Aims*: This study aimed to explore the use of a new set of tools from COSMIN to appraise

8    studies on outcome measures available to speech and language therapists (SLTs).

9    *Methods*: We used the COSMIN tools to appraise seven studies and a user manual that reported

10   the measurement properties of the Focus on the Outcomes of Children Under Six (FOCUS), a

11   validated measure of preschoolers' communicative participation that is used in various contexts

12   around the world.

13   *Results*: Using COSMIN guidelines, the FOCUS was categorized as a "Category A" tool because

14   there was a sufficient level of evidence to support its content validity and internal consistency.

15   According to the COSMIN guidelines, this means that the FOCUS can be recommended for

16   clinical use. The quality of evidence supporting measurement properties of the FOCUS received

17   a rating of 'moderate', meaning users can have moderate confidence in its measurement

18   properties. Since these ratings from the COSMIN tools may be unclear to users of the FOCUS,

19   we have provided more specific recommendations.

20   *Conclusions & Implications*: The COSMIN tools offer detailed standards to support the appraisal

21   of outcome measures available to SLTs. However, several limitations were observed, and

22   recommendations to support the application of the COSMIN tools are provided.

23

24   **Declaration of Interest**: The authors declare no conflicts of interest.

25

**What is already known on this subject?**

Collecting outcome data is essential to ensuring speech and language therapy is effective. Until

the development of Consensus-based Standards for the selection of health Measurement

INstruments (COSMIN) there was a lack of standards in the way the measurement properties of

outcome measures were appraised.

**What this study adds?**

This paper used the Focus on the Outcomes of Children Under Six (FOCUS), a measure of

preschoolers' communicative participation outcomes in speech and language therapy, as a case

example to illustrate the applications of the COSMIN tools. In doing so the strengths and

limitations of the current COSMIN tools in appraising the quality of outcome measure

instruments are emphasized.

**Clinical implications of this study?**

The COSMIN tools offer a step-by-step, standardized approach to appraising various

measurement properties in outcome instruments. Due to existing limitations of the COSMIN

tools, appraisal should provide clear and specific recommendations so users of outcome

measures (e.g., SLTs, researchers) can identify the appropriate uses of each instrument.

49 **Introduction**

50       Outcome measures are important tools for assessing the impact of a healthcare system

51 (Agency for Health Research and Quality 2011, Donabedian 1988). Across the globe, speech

52 language therapists (SLTs) are encouraged by their professional organizations to use outcome

53 measures (Mullen and Schooling 2010, Royal College of Speech & Language Therapists 2020,

54 Speech-Language & Audiology Canada 2010). Amongst many benefits, data collected using

55 outcome measures allow for the evaluation of clinical effectiveness, inform quality improvement

56 efforts, and support best practices (Royal College of Speech & Language Therapists 2020).

57 Moreover, when used in large health systems, data collected using valid and reliable outcome

58 measures can generate evidence to inform decisions about services (e.g., service type, length, and

59 intensity). For SLTs, outcome measures can be used to gather objective data on clients' skills

60 and progress – which can be used to guide clinical decisions (Garland *et al.* 2003). Patient-

61 reported outcome measures provide clients and families with opportunities to express their

62 perspectives, values, and preferences about their own care, improving SLTs' accountability to

63 their clients (Ronen *et al.* 2000).

64       To realize the many benefits associated with outcome measures, it is imperative to select

65 tools that have appropriate measurement properties (Enderby and John 2015, 2020, Speech-

66 Language & Audiology Canada 2012, Threats 2013, World Health Organization 2001). Despite

67 some graduate training to support understanding of psychometrics, in practice SLTs report a lack

68 of confidence (Kerr *et al.* 2003), time (Kerr *et al.* 2003), and resources (e.g., access to literature

69 Vallino-Napoli and Reilly, 2004) to evaluate the properties of the outcome measures they use.

70 These barriers may explain why measurement properties of instruments were not a major factor

71 influencing SLTs' choice of instrument (Betz *et al.* 2013) and that "misuses" of measurement

72  instruments were frequently reported (Kerr *et al.* 2003). One way to support SLTs in choosing

73  appropriate instruments is to appraise existing outcome measures systematically and critically.

74       The Consensus-based Standards for the selection of health Measurement INstruments

75  (COSMIN) offers an internationally agreed-upon taxonomy for evaluating the psychometric

76  properties of health-related patient reported outcome measures (Mokkink, Terwee, Patrick, *et al.*

77  2010). According to COSMIN, an outcome measure can be evaluated based on (1) the methods

78  used in *tool development* and (2) psychometric properties (*validity*, *reliability*, and

79  *responsiveness*) (Barten *et al.* 2012, Lambert and Hawkins 2004, Mokkink, Terwee, Patrick, *et al.*

80  2010). Additionally, COSMIN acknowledges the importance of two additional properties:

81  *interpretability* and *feasibility*. Interpretability refers to the ease of deriving meaning from an

82  instrument's scores, and feasibility refers to the ease with which an instrument is adopted into its

83  intended context (Mokkink *et al.* 2018, Mokkink, Terwee, Patrick, *et al.* 2010). At the moment,

84  *interpretability* and *feasibility* are not formally assessed with the COSMIN tools, but rather are

85  viewed as important considerations for the practical use of an outcome measure (Mokkink,

86  Terwee, Patrick, *et al.* 2010).

87       To support the evaluation of patient-reported outcome measures, COSMIN offers a

88  validated Risk of Bias checklist (Mokkink *et al.* 2018) and a user manual with step-by-step guide

89  to support instrument appraisals (COnsensus-based Standards for the selection of health

90  Measurement INstruments 2020, Prinsen *et al.* 2018, Terwee *et al.* 2018). The Risk of Bias

91  checklist was developed based on a literature review of the measurement properties of health-

92  related measurement instruments and the consensus of a panel of 57 experts involved in a Delphi

93  study (Mokkink *et al.* 2018, Mokkink, Terwee, Knol, *et al.* 2010). Although the COSMIN tools

94  were originally developed and validated to appraise patient-reported outcome measures, it has

95  been argued that their criteria  are also applicable to the evaluation of non-patient reported

96    outcome measures (Tate 2019). Since becoming available, these COSMIN tools have been used

97    to evaluate a range of patient/parent/clinician-reported outcome measures within healthcare (e.g.,

98    Bull et al., 2019; Howell et al., 2020; Williams and Beovich, 2020), and were found to be some

99    of the more carefully developed and comprehensive tools to appraise outcome measures (Tate

100   2019). Additionally, some professional organizations have begun to use COSMIN tools to

101   recommend instruments that met standards for clinical use (England *et al.* 2019, Pick *et al.* 2020).

102        The purpose of this paper is to illustrate use of the COSMIN tools (i.e., the Risk of Bias

103   checklist) (Mokkink *et al.* 2018) and the criteria for good measurement properties (Prinsen *et al.*

104   2018, Terwee *et al.* 2018) to appraise an outcome measure in speech-language therapy.

105   Importantly, this paper illustrates *how* appraisal results can be considered to draw clinically

106   meaningful recommendations regarding the use of existing outcome measures. We argue that

107   existing outcome measures should be considered on an instrument-by-instrument basis for three

108   reasons. First, most existing outcome measures are found not to meet all standards of the

109   COSMIN tools, which is not surprising given that most tools were developed prior to the

110   COSMIN standards (Bull *et al.* 2019, Howell *et al.* 2020, Williams and Beovich 2020). Second,

111   not all measurement properties are equally important for all clinical decisions, so the properties

112   important to individual clinical decisions or purposes should be considered in context, on a

113   measure-by-measure basis rather than categorizing tools as 'good' versus 'bad' (Bull *et al.* 2019,

114   Messick 1993, 1995). Third, in practice, clinicians are limited by the resources available to them

115   (i.e., instruments available in their clinic/district). Therefore, considering what each outcome

116   measure can and cannot do is perhaps more practical and appropriate than identifying one "best"

117   tool.

118        To contextualize the considerations when applying the COSMIN tools, we have included

119   an evaluation of an outcome measure that is currently implemented in at least one large clinical

120    health system. The Focus on the Outcomes of Children Under Six (FOCUS) (Thomas-Stonell *et*

121    *al.* 2015) is one of a handful of validated tools explicitly designed to measure *outcomes* for

122    preschoolers receiving SLT interventions (Thomas-Stonell *et al.* 2010). Furthermore, the

123    FOCUS is the only validated tool available to assess how preschoolers use their communication

124    to participate in real-world situations (Cunningham *et al.* 2017) – an outcome that has been

125    identified as important and meaningful to families (Lindsay and Dockrell 2004, Roulstone *et al.*

126    2013).

127          With the FOCUS as a case example, the goal of this paper is to illustrate how the

128    COSMIN tools can be used as a guide to identify the strengths and limitations that are associated

129    with any outcome measure. We intend for this paper to serve as a support for researchers and

130    SLTs in selecting tools that are both psychometrically strong and meaningful for practice. The

131    paper will also have implications for test developers who will want to understand a new standard

132    for evaluating outcome measures. This work involves the secondary analysis of data collected

133    during a recent scoping review of the literature related to the FOCUS (Cunningham *et al.* 2020).

134    More specifically, we used findings from this review to identify studies that reported

135    psychometric properties of the FOCUS, which we appraised using the COSMIN tools.

136    **Methods**

137    *Search strategy and inclusion criteria*

138          Cunningham et al. identified 25 publications that reported on either the development or

139    application of the FOCUS (Cunningham *et al.* 2020). In the current study, we reviewed these 25

140    publications to identify those that reported psychometric properties of the FOCUS using the

141    following inclusion criteria: (i) the article was peer-reviewed and about the English-language

142    version of the FOCUS or its derivative (i.e., the FOCUS-34); (ii) the study evaluated the

143    psychometric properties of the FOCUS; and (iii) the article was published in English. Although

144    not peer reviewed, the published FOCUS user manual was also included in order to complete a

145    comprehensive appraisal (Thomas-Stonell *et al.* 2020). This user manual presented information

146    for both the original and shortened FOCUS tools (i.e., FOCUS-34), and was created by drawing

147    upon peer-reviewed research.

148    *Appraisal of psychometric properties.*

149        An extraction spreadsheet (see Appendix 1) was developed to record the following data

150    from each identified publication: (1) FOCUS version, (2) measurement properties investigated,

151    and (3) study methodology and results. Based on the psychometric properties investigated in

152    each publication, the relevant portions of the COSMIN Risk of Bias checklist (Mokkink *et al.*

153    2018) and the criteria for good measurement properties were completed (i.e., the reliability

154    portion for studies that investigated reliability) (Prinsen *et al.* 2018, Terwee *et al.* 2018).

155    Descriptions of each measurement property as well as explanations of how each applies to the

156    FOCUS are presented in table 1.

157                              < Table 1 Here >

158        ***Appraisal of a validation study's methodology***. The COSMIN Risk of Bias checklist has

159    specific sections dedicated to evaluating the *methodology* of studies conducted to demonstrate

160    measurement properties of an outcome measure. Under each section (i.e., for each specific

161    measurement property), the checklist provides a list of items to evaluate a study's quality. For

162    example, the section on content validity contains 31 items concerning the appropriateness of data

163    collection methods, participant sample, sample size, and data analysis approach. For each item, a

164    study's methodology is rated as *very good*, *adequate*, *doubtful* or *inadequate* based on well-

165    defined standards (these ratings corresponding to *excellent*, *good*, *fair*, *poor* on the original rating

166    scale (Mokkink *et al.* 2018)). For example, the standards for sample sizes are: ≥50 (for *very good*

167    rating), ≥30 (*adequate*), <30 (*doubtful*), unclear sample size (*inadequate*). Across all items, the

168 lowest rating is selected as the ratings for a study's methodologies (i.e., the worst score counts)

169 (Mokkink *et al.* 2018). In cases where multiple studies evaluated the same measurement property

170 and received different quality ratings, the higher quality rating was taken as the overall rating.

171 We reasoned that measurement properties are subject to continuous evaluation; thus, if multiple

172 studies were conducted over several years, the study with the most rigorous design should be

173 used. Data extraction and quality rating were first completed by E.K. and reviewed by NT-S with

174 93% agreement. Due to the "worst score counts" rule, no disagreements resulted in changes in

175 the rating of any psychometric property. All disagreements were resolved through discussion

176 with BJC and PR.

177     ***Appraisal for validation study's results***. Apart from evaluating the methodology of

178 validation studies, the results of each primary study were considered using the COSMIN criteria

179 for good measurement properties (Prinsen *et al.* 2018, Terwee *et al.* 2018). This provides a

180 quality rating based on the results reported in validation studies regarding a psychometric

181 property of an outcome measure, and can range from "+" sufficient; "-" insufficient; "?"

182 indeterminate; "±" inconsistent. For example, for reliability testing, a study would receive a "+"

183 if it found an ICC or weighted Kappa value $\geq 0.70$. In contrast, a "–" rating would be assigned if

184 the reported ICC or weighted Kappa value was <0.70. A "?" would be assigned if no ICC or

185 weighted Kappa was reported. If more than one study was conducted for reliability but had

186 mixed findings, an "±" inconsistent rating was assigned. For the evaluation of a study's results,

187 E.K. and NT-S had 100% agreement in their ratings.

188     Together, these COSMIN tools provide standards to appraise study methodologies and

189 reported results (Mokkink *et al.* 2018). There are, however, some limitations to the COSMIN

190 rating scales. When validating the COSMIN Risk of Bias rating, the COSMIN developers noted

191 that "*a study often received a 'fair' quality rating (i.e., doubtful on the new rating scale) only*

192 *because it was not reported how missing items were handled. It was argued that this would not*

193 *necessarily lead to biased results of the study*" (Mokkink et al., 2018, p.2). This is why, when a

194 study receives *doubtful/inadequate* rating on Risk of Bias or an *indeterminate* rating on the

195 criteria for good measurement properties (Prinsen *et al.* 2018, Terwee *et al.* 2018), the reasons

196 behind the rating should be scrutinized before assigning an overall quality rating for the

197 measurement tool. To this end, COSMIN also provides a rating scale (*High*, *Moderate*, *Low*,

198 *Very Low*) and key factors to consider when indicating the overall quality of a measurement tool.

199 The four key factors to consider include: high risk of bias in study methodology and reporting,

200 inconsistent findings across studies, imprecision (referring to a small sample size) and

201 indirectness (referring to validation studies completed in a population dissimilar from the

202 intended users of the instrument).

203　　　Lastly, to improve transparency, COSMIN recommends categorizing outcome measures

204 into the following categories: (A) instrument with evidence for sufficient content validity AND

205 at least low-quality evidence for sufficient internal consistency; (B) instrument categorized not in

206 A or C; and (C) instrument with high quality evidence for an insufficient measurement property.

207 According to COSMIN, outcome measures in Category A can be recommended for use and their

208 results can be trusted; those in category B can be recommended provisionally, subject to further

209 evidence being provided; and category C tools should not be recommended.

210 **Results**

211　　　Full-text screening identified 7 articles that met the inclusion criteria. These articles and

212 the FOCUS user manual were included in the appraisal (reasons for studies being excluded are

213 shown in figure 1). The studies included described the methodology used to develop FOCUS

214 items (Thomas-Stonell *et al.* 2009), content validity (Oddson *et al.* 2019, Thomas-Stonell *et al.*

215 2010), construct validity (Thomas-Stonell *et al.* 2010, Washington, Thomas-Stonell, *et al.* 2013),

216   internal-consistency reliability (Oddson *et al.* 2019, Thomas-Stonell *et al.* 2010), inter-rater

217   reliability (Oddson *et al.* 2013, Thomas-Stonell *et al.* 2010, 2013, Washington, Oddson, *et al.*

218   2013), test-retest reliability (Thomas-Stonell *et al.* 2010, Washington, Oddson, *et al.* 2013),

219   responsiveness (Thomas-Stonell *et al.* 2013), and interpretability (Oddson *et al.* 2019, Thomas-

220   Stonell *et al.* 2013) of the FOCUS. It should be noted that the FOCUS-34 is a streamlined

221   version of the original 50-item FOCUS, as 16 items were removed based on empirical findings

222   from item response analysis (Oddson *et al.* 2019). For readability, measurement properties that

223   apply to both the original FOCUS and the FOCUS-34 are described as *FOCUS tools*.

224                               < Insert Figure 1 Here >

225        The COSMIN quality ratings for the FOCUS validation studies are presented in table 2.

226   Considering all available evidence related to the measurement properties of the FOCUS, we

227   rated the quality of evidence as *Moderate* and categorized the FOCUS tools as Category A due to

228   sufficient content validity and internal consistency. According to COSMIN, this means that the

229   FOCUS tools can be recommended for use, and that we are *moderately confident* that the

230   FOCUS provides an estimate close to what has been stated in the literature (i.e., the reported

231   measurement properties). The major considerations that led to this overall rating are described in

232   the sections that follow. A detailed rationale behind each quality rating presented in table 2 can

233   be found in appendix 1. Item-by-item scoring of the COSMIN tools is available from the authors

234   upon request. We acknowledge that the overall *Moderate* rating and Category A nomenclature

235   are not very informative, therefore, we provided the following sections to describe the clinical

236   implications from the appraisal findings.

237        ***Tool Development and Content Validity*** (Oddson *et al.* 2019, Thomas-Stonell *et al.* 2009,

238   2010): Current findings suggest that the FOCUS measures communicative participation

239   outcomes that are meaningful and important to both parents and SLTs. Each of the FOCUS items

240    was found to be clear and relevant to users (the development and testing of the FOCUS items

241    involved 349 parents and SLTs). Through a combination of quantitative and qualitative analysis,

242    the items on the FOCUS-34 were also found to provide a comprehensive measure of

243    communicative participation outcomes. These studies received positive ratings because a clear

244    description was provided of the aim of the FOCUS, the high ecological validity in the item

245    generation, selection and reduction process that involved parents, SLTs and statisticians. For

246    SLTs, since the FOCUS was validated in various real-world clinical settings that serve a range of

247    clinical populations, it provides a consistent tool to measure gains in preschoolers'

248    communicative participation skills during speech-language interventions.

249            *Validity* (Oddson *et al.* 2019, Thomas-Stonell *et al.* 2010, 2013, Washington, Thomas-

250    Stonell, *et al.* 2013): FOCUS scores were shown to correlate moderately with existing

251    instruments that measure related, but dissimilar constructs (e.g., the Pediatric Quality of Life

252    Inventory (Varni 1998), the communication domains on the Ages and Stages Questionnaire –

253    Social/ Emotional (ASQ-SE) (Squires *et al.* 2003), Communication and Socialization domains of

254    the Vineland Adaptive Behavior Scales (VABS) (Sparrow *et al.* 2005)). Meanwhile, FOCUS

255    scores did not correlate with domains that are not related to communication. These studies on

256    construct validity received an overall *very good* rating on methodology and *sufficient* rating on

257    results. For SLTs, these findings clarify what is being measured by the FOCUS, namely an

258    aspect of communication that relates to how children use communication skills to participate in

259    everyday situations. Additionally, there is a very high correlation between scores on the FOCUS-

260    34 and the original FOCUS in the criterion validity study (Oddson *et al.* 2019), which suggests

261    that the FOCUS-34 sufficiently reflects the original tool and can provide a more efficient option

262    for data collection for those who need or want it.

263        ***Reliability*** (Oddson *et al.* 2013, 2019, Thomas-Stonell *et al.* 2010, Washington, Oddson,

264    *et al.* 2013): Two studies reported on internal consistency (Oddson *et al.* 2019, Thomas-Stonell

265    *et al.* 2010) and received an overall rating of *very good* for methodology and *sufficient* for results.

266    Studies that explored inter-rater reliability received an overall *adequate* rating for methodology

267    and *sufficient* rating for results. There was a moderate level of inter-rater reliability between

268    SLTs and parents or amongst SLTs, which suggests that even when completed by different

269    individuals, scores on the FOCUS reliably reflect preschoolers' communicative participation

270    skills. Thus, it is not necessary for both parents and SLTs to complete the FOCUS in order to

271    capture change. With regards to test-retest reliability, the current *doubtful* and *indeterminate*

272    ratings were due to the fact that Pearson's correlation ($r = 0.96$) (Washington, Oddson, *et al.* 2013)

273    were reported instead of ICC values in the FOCUS validation studies. While the quality ratings

274    were limited by the reported statistics, we do not believe it t should not limit use of the FOCUS

275    in clinical practice.

276        ***Responsiveness*** (Thomas-Stonell *et al.* 2013): There is no single agreed-upon approach

277    for measuring responsiveness (i.e., an outcome measure's ability to detect change) (Thomas-

278    Stonell *et al.* 2007) and COSMIN offers a range of checklists to assess responsiveness of an

279    outcome measure. The responsiveness of the FOCUS was evaluated in two ways. First, change

280    scores on the FOCUS were compared to the change scores measured by three established

281    measures of speech, intelligibility and language (i.e., Children's Speech Intelligibility Measure

282    (Wilcox and Morris 1999), Percent Consonant Correct-Revised (Schriberg *et al.* 1997) and

283    Developmental Sentence Scoring (Lee and Canter 1971)) (Thomas-Stonell *et al.* 2013). There

284    was a fair level of agreement between measures when a minimally clinically important

285    difference was observed. A fair but not excellent level of agreement is to be expected, since the

286    FOCUS and these comparator measures do not measure the same construct. This study received

287     a *very good* rating for methodology and *sufficient* rating for results. Second, responsiveness of

288     the FOCUS was demonstrated when preschoolers receiving interventions showed more change

289     than a group of children on a waitlist ($M =18.2$ and $M = 5.87$ points respectively, and that the

290     average change scores in the waitlist group was lower than the 16-point cut-off scores to be

291     considered minimally clinically significant change) (Thomas-Stonell *et al.* 2013). However,

292     because the study was observational rather than experimental, and was conducted within a

293     practice context, this finding was limited by the unequal intervention (90 days) and waitlist (60

294     days) intervals. This resulted in a *doubtful* rating for the study's methodology. For the FOCUS-

295     34 tool, the change scores (i.e., pre-to-post intervention) highly correlate with those from the

296     original FOCUS tool.

297        For SLTs, these findings suggest that the FOCUS demonstrates comparable

298     responsiveness with speech, intelligibility and language outcome measures that are commonly

299     used in practice. One area of constraint related to responsiveness is that with the published

300     evidence, the possibility that changes observed on the FOCUS were due to natural development

301     or some other factors cannot be fully ruled out. However, SLTs using the recommended criterion

302     to interpret when a minimally clinically significant change has occurred will minimize

303     contributions from natural development and random error. In the absence of a control group, we

304     caution SLTs against attributing change on the FOCUS solely to specific treatment effects, as

305     this is difficult to determine given the many factors that can affect children's development at any

306     given time (e.g., growth spurt, change in language learning environment).

307        ***Feasibility & Interpretability*** (Thomas-Stonell *et al.* 2020): These two properties are not

308     formally evaluated by the COSMIN tools, so here we summarize the major findings related to

309     the FOCUS for these two properties. The streamlined FOCUS-34 provides a reliable and

310     efficient option for data collection, which can be completed by parents or SLTs within 10

311  minutes. The FOCUS tools offer criterion scores to support SLTs in interpreting change in

312  children's communicative participation during intervention. The criterion score had 95%

313  agreement between parents' and SLTs' judgements of whether a clinically important change had

314  occurred. This criterion score allows SLTs and researchers to determine whether meaningful

315  change occurred during an intervention period, and theoretically this can be done without a

316  control group, making it particularly useful for both research and practice.

317                          <Insert Table 2 Here>

318  **Discussion:**

319        The purpose of this paper was to illustrate the use of relatively new tools from the

320  COSMIN (Mokkink *et al.* 2018) to guide the appraisal of outcome measurement tools. As such,

321  the discussion is focused on the benefits and limitations of the COSMIN.

322        The COSMIN tools (2018) are comprehensive and offer standards to appraise patient-

323  reported outcome measures s. For trainees and SLTs, the standards provide an objective way of

324  appraising measurement properties of outcome measures. We reiterate that the appraisal should

325  be conducted on an instrument-by-instrument basis and guided by a clearly articulated clinical or

326  research question(s). For tool developers, COSMIN provides a standard to improve the quality of

327  reporting for the development and validation of outcome measurement tools. In fact, two authors

328  of this paper (NS-T and PR) were involved in the development and validation of the FOCUS

329  tools, and this checklist has helped identify additional details that could be included in future

330  editions of the FOCUS user's manual to continue to support clinical practice and research.

331        We also observed important limitations in our efforts to apply the COSMIN tools. The

332  first is the time needed to complete the appraisal. It took over 25 hours for our team of four

333  authors with graduate-level training in tool development to complete the evaluation of the

334  FOCUS tools (not including the time to become acquainted with the COSMIN tools). Clinicians

335     focused on providing quality care to clients are unlikely to have the time or academic

336     background to complete this type of detailed evaluation (e.g., the knowledge to evaluate

337     statistical analysis and results such as item-response analysis). Thus, the onus of evaluating

338     measurement properties may necessarily fall to interested researchers, professional colleges and

339     tool developers.

340         Second, the COSMIN tool has not yet developed a rating scale to evaluate *interpretability*

341     or *feasibility*, but we believe these are among the most important clinical considerations for SLTs;

342     they are interested in understanding whether observed changes are clinically meaningful, and

343     whether a tool can be easily adopted into practice. When using COSMIN to appraise outcome

344     measurement tools, these two properties should not be overlooked simply due to a lack of clear

345     appraisal standards. Until a rating scale is available on COSMIN, we recommend referring to the

346     detailed data extraction matrix that is available in the COSMIN user manual (Prinsen *et al.* 2018,

347     Terwee *et al.* 2018) to identify information related to interpretability and feasibility. We also

348     recommend using other tools to supplement appraisals in these areas, for example, the criteria

349     from the *Acceptability and Utility* checklist from the Allied Health Professions (AHP) Outcome

350     Measures UK Working Group (Allied Health Professions (AHP) Outcome Measures UK

351     Working Group, 2019) and the *Interpretability* and *Burden* tool from the Scientific Advisory

352     Committee of the Medical Outcomes Trust (Lohr 2002).

353         Thirdly, we emphasize the need for any appraisal completed using the COSMIN tools to

354     consider the practical implications of appraisal findings (i.e., making clear recommendations

355     regarding tool use). One important reason for making a clear recommendation statement is

356     concern about the categories on the COSMIN risk of bias checklist (i.e., the *very good*, *adequate*,

357     *doubtful* and *inadequate* scale) and the criteria for determining good measurement properties (i.e.,

358     the + sufficient, ? indeterminate , - insufficient). Using COSMIN, each measurement property

359  receives ratings that reflect the design and reporting of validation studies, and not the

360  measurement property of the tool itself. Clinicians, policymakers and other stakeholder groups

361  who are unfamiliar with the COSMIN tools may take these ratings to mean that a tool is "very

362  good", "adequate", "doubtful" or "inadequate". Current knowledge on measurement properties

363  suggests that users consider the *purpose for measurement* (i.e., What is the clinical/research

364  question?*)*. A measurement tool is 'good' or 'bad' for specific usesl, but  tool should not be

365  viewed as categorically good versus bad (Bull *et al.* 2019, Messick 1993, 1995). Another reason

366  for making clear recommendations is the fact that systematic reviews have reported many

367  existing outcome measurement tools would not be considered to be of adequate/sufficient quality

368  based on appraisals done using the COSMIN tools (Bull *et al.* 2019, Howell *et al.* 2020,

369  Williams and Beovich 2020). This is a commonly reported limitation in studies that

370  retrospectively applied the COSMIN tools to evaluate a measurement tool that was developed

371  prior to COSMIN being published; it reflects a lack of standards in the reporting of measurement

372  properties as well as an evolving understanding of best-practice in tool development and

373  validation (Bull *et al.* 2019, Van Tiggelen *et al.* 2020, Williams and Beovich 2020). Providing

374  clear recommendations will help users interpret appraisal findings accurately, and understand the

375  appropriate use of existing outcome measurement tools. The results reported in this paper serve

376  as a case example for how clear recommendations can (and should) be made on an instrument-

377  by-instrument basis, depending on the purpose of the measure and the question(s) to be answered.

378          Another limitation of COSMIN relates to the scope of application. The COSMIN tools

379  were originally developed and validated to appraise *patient-reported* outcome measures. While it

380  has been argued that the criteria in the COSMIN tools are applicable to evaluate non-patient

381  reported outcome measures (Tate 2019), it is possible that more criteria should be considered

382  when appraising non-patient reported outcome measures. Recent work is expanding the scope of

383     the COSMIN tools for the appraisal of clinician-reported, performance-based and laboratory-

384     based outcome measure instruments (Mokkink *et al.* 2020). As such, when more comprehensive,

385     validated appraisal tools become available, the work described here will be updated.

386         A future direction of our work is to appraise multiple functional outcome measures used

387     by SLTs, particularly those used with young children. These appraisals will allow us to identify

388     the strengths and limitations, and the specific purposes, of existing outcome measures, and the

389     appropriate uses of each of the available measures for SLTs. Recommendations will be

390     developed based on these appraisals, and will be available as an online resource for  SLTs.

391     **Conclusion:**

392         The study illustrates the use, and limitations, of the COSMIN tools (Mokkink *et al.* 2018,

393     Prinsen *et al.* 2018, Terwee *et al.* 2018), which were designed to appraise outcome measures

394     systematically. The COSMIN tools provide an up-to-date, comprehensive list of factors to

395     consider in psychometric appraisals, but due to an evolving understanding of psychometric

396     properties and reporting standards, many existing clinical tools (i.e., those developed prior to the

397     COSMIN tools) may receive doubtful/indeterminate ratings on COSMIN. Appraisal of all

398     existing outcome measurement tools should consider carefully the reasons behind quality ratings

399     and how these may impact clinical decisions. This paper demonstrates how measurement

400     properties should be considered in conjunction with clinical decisions to be made based on using

401     the outcome measurement instrument(s). Lastly, for researchers and tool developers, this paper

402     introduces a newly available tool that can be used to guide the development and reporting of

403     outcome measurement instrument(s). We believe this study will be a useful reference for SLTs,

404     researchers, and developers in appraising, choosing and creating appropriate outcome

405     measurement tools.

406

407                                                        References

408    AGENCY FOR HEALTH RESEARCH AND QUALITY, 2011, Types of health care quality measures

409        [online]. Available: https://www.ahrq.gov/talkingquality/measures/types.html [Accessed 8

410        January 2020].

411    ALLIED HEALTH PROFESSIONS (AHP) OUTCOME MEASURES UK WORKING GROUP, 2019, Key

412        questions to ask when selecting outcome measures: a checklist for allied health

413        professionals.

414    BARTEN, J.A., PISTERS, M.F., HUISMAN, P.A., TAKKEN, T., and VEENHOF, C., 2012,

415        Measurement properties of patient-specific instruments measuring physical function.

416        *Journal of Clinical Epidemiology*, **65**, 590–601.

417        https://doi.org/10.1016/j.jclinepi.2011.12.005.

418    BETZ, S.K., EICKHOFF, J.R., and SULLIVAN, S.F., 2013, Factors influencing the selection of

419        standardized tests for the diagnosis of specific language impairment. *Language, Speech, and*

420        *Hearing Services in Schools*, **44**, 133–146. https://doi.org/10.1044/0161-1461(2012/12-

421        0093).

422    BULL, C., BYRNES, J., HETTIARACHCHI, R., and DOWNES, M., 2019, A systematic review of the

423        validity and reliability of patient-reported experience measures. *Health Services Research*,

424        **54**, 1023–1035. https://doi.org/10.1111/1475-6773.13187.

425    CONSENSUS-BASED STANDARDS FOR THE SELECTION OF HEALTH MEASUREMENT INSTRUMENTS,

426        2020, COSMIN Tools [online]. Available: https://www.cosmin.nl/cosmin-tools/ [Accessed

427        2 November 2020].

428    CUNNINGHAM, B.J., THOMAS-STONELL, N.L., and ROSENBAUM, P., 2020, Assessing

429        communicative participation in preschool children with the Focus on the Outcomes of

430        Communication Under Six: a scoping review. *Developmental Medicine & Child Neurology*.

431       https://doi.org/10.1111/dmcn.14665.

432   CUNNINGHAM, B.J., WASHINGTON, K.N., BINNS, A., ROLFE, K., ROBERTSON, B., and

433       ROSENBAUM, P., 2017, Current methods of evaluating speech-language outcomes for

434       preschoolers with communication disorders: A scoping review using the ICF-CY. *Journal*

435       *of Speech, Language, and Hearing Research*, **60**, 447–464. https://doi.org/10.1044/2016.

436   DONABEDIAN, A., 1988, The quality of care: How can it be assessed? *JAMA*, **260**, 1743–1748.

437       https://doi.org/10.1001/jama.1988.03410120089033.

438   ENDERBY, P. and JOHN, A., 2015, *Therapy outcome measures for rehabilitation professionals*.

439       Third Edit. (Guildford: J&R Press Ltd).

440   ENDERBY, P. and JOHN, A., 2020, *Therapy Outcome Measures theoretical underpinning and case*

441       *studies* (Havant: J&R Press Ltd).

442   ENGLAND, B.R., TIONG, B.K., BERGMAN, M.J., CURTIS, J.R., KAZI, S., MIKULS, T.R., O'DELL,

443       J.R., RANGANATH, V.K., LIMANNI, A., SUTER, L.G., and MICHAUD, K., 2019, 2019 Update

444       of the American College of Rheumatology recommended rheumatoid arthritis disease

445       activity measures. *Arthritis Care and Research*, **71**, 1540–1555.

446       https://doi.org/10.1002/acr.24042.

447   GARLAND, A.E., KRUSE, M., and AARONS, G.A., 2003, Clinicians and outcome measurement:

448       What's the use? *Journal of Behavioral Health Services and Research*, **30**, 393–405.

449       https://doi.org/10.1007/BF02287427.

450   HOWELL, M., BRADSHAW, J., and LANGDON, P.E., 2020, A Systematic Review of Behaviour-

451       Related Outcome Assessments for Children on the Autism Spectrum with Intellectual

452       Disabilities in Education Settings. *Review Journal of Autism and Developmental Disorders*.

453       https://doi.org/10.1007/s40489-020-00205-y.

454   KERR, M.A., GUILDFORD, S., and BIRD, E.K.., 2003, Standardized language test use: A Canadian

455    survey utilisation. *Journal of Speech-Language Pathology and Audiology*, **27**, 10–28.

456  LAMBERT, M.J. and HAWKINS, E.J., 2004, Measuring outcome in professional practice:

457    Considerations in selecting and using brief outcome instruments. *Professional Psychology:*

458    *Research and Practice*, **35**, 492–499. https://doi.org/10.1037/0735-7028.35.5.492.

459  LEE, L. and CANTER, S., 1971, Developmental sentence scoring: a clinical procedure for

460    estimating syntactic development in children's spontaneous speech. *J Speech Lang Hear*

461    *Disord*, **36**, 315–40.

462  LINDSAY, G. and DOCKRELL, J.E., 2004, Whose job is it? Parents' concerns about the needs of

463    their children with language problems. *Journal of Special Education*, **37**, 225–235.

464    https://doi.org/10.1177/00224669040370040201.

465  LOHR, K.N., 2002, Assessing health status and quality-of-life instruments: Attributes and review

466    criteria. *Quality of Life Research*, **11**, 193–205. https://doi.org/10.1023/A:1015291021312.

467  MESSICK, S., 1993, Foundations of validity: Meaning and consequences of psychological

468    assessment. *Educational Testing Service Research Report Series*, **2**.

469  MESSICK, S., 1995, Standards of validity and the validity of standards in performance assessment.

470    *Educational Measurement: Issues and Practice*, **14**, 5–8.

471  MOKKINK, L., BOERS, M., VLEUTEN, C. VAN DER, BOUTER, L., ALONSO, J., PATRICK, D., VET, H.

472    DE, and TERWEE, C., 2020, COSMIN Risk of Bias tool to assess the quality of studies on

473    reliability or measurement error of outcome measurement instruments: a Delphi study. *BMC*

474    *Medical Research Methodology*, **1**, 1–13. https://doi.org/10.21203/rs.3.rs-40864/v1.

475  MOKKINK, L.B., TERWEE, C.B., KNOL, D.L., STRATFORD, P.W., ALONSO, J., PATRICK, D.L.,

476    BOUTER, L.M., and DE VET, H.C.W.W., 2010, The COSMIN checklist for evaluating the

477    methodological quality of studies on measurement properties: A clarification of its content.

478    *BMC Medical Research Methodology*, **10**. https://doi.org/10.1016/j.jclinepi.2010.02.006.

479    MOKKINK, L.B., TERWEE, C.B., PATRICK, D.L., ALONSO, J., STRATFORD, P.W., KNOL, D.L.,

480        BOUTER, L.M., and DE VET, H.C.W., 2010, The COSMIN study reached international

481        consensus on taxonomy, terminology, and definitions of measurement properties for health-

482        related patient-reported outcomes. *Journal of Clinical Epidemiology*, **63**, 737–745.

483        https://doi.org/10.1016/j.jclinepi.2010.02.006.

484    MOKKINK, L.B., DE VET, H.C.W., PRINSEN, C.A.C., PATRICK, D.L., ALONSO, J., BOUTER, L.M.,

485        and TERWEE, C.B., 2018, COSMIN Risk of Bias checklist for systematic reviews of Patient-

486        Reported Outcome Measures. *Quality of Life Research*, **27**, 1171–1179.

487        https://doi.org/10.1007/s11136-017-1765-4.

488    MULLEN, R. and SCHOOLING, T., 2010, The national outcomes measurement system for pediatric

489        speech-language pathology. *Language Speech and Hearing Services in Schools*, **41**, 44.

490        https://doi.org/10.1044/0161-1461(2009/08-0051).

491    ODDSON, B., THOMAS-STONELL, N.L., ROBERTSON, B., and ROSENBAUM, P., 2019, Validity of a

492        streamlined version of the Focus on the Outcomes of Communication Under Six: Process

493        and outcome. *Child: Care, Health and Development*, 600–605.

494        https://doi.org/10.1111/cch.12669.

495    ODDSON, B., WASHINGTON, K.N., ROBERTSON, B., THOMAS-STONELL, N., and ROSENBAUM, P.,

496        2013, Inter-rater reliability of clinicians' ratings of preschool children using the FOCUS©:

497        Focus on the outcomes of communication under six. *Canadian Journal of Speech-Language*

498        *Pathology and Audiology*, **37**, 170–174.

499    PICK, S., ANDERSON, D.G., ASADI-POOYA, A.A., ASADI-POOYA, A.A., AYBEK, S., BASLET, G.,

500        BLOEM, B.R., NICHOLSON, T.R., BROWN, R.J., CARSON, A.J., CHALDER, T., DAMIANOVA,

501        M., DAVID, A.S., EDWARDS, M.J., EPSTEIN, S.A., ESPAY, A.J., GARCIN, B., GOLDSTEIN, L.H.,

502        HALLETT, M., JANKOVIC, J., JOYCE, E.M., KANAAN, R.A., KEYNEJAD, R.C., KOZLOWSKA, K.,

503    LAFAVER, K., CURT LAFRANCE, W., LANG, A.E., LEHN, A., LIDSTONE, S., MAURER, C.W.,

504    MILDON, B., MORGANTE, F., MYERS, L., NICHOLSON, C., NIELSEN, G., PEREZ, D.L.,

505    POPKIROV, S., REUBER, M., ROMMELFANGER, K.S., SCHWINGENSHUH, P., SERRANOVA, T.,

506    SHOTBOLT, P., STEBBINS, G.T., STONE, J., TIJSSEN, M.A.J., and TINAZZI, M., 2020, Outcome

507    measurement in functional neurological disorder: A systematic review and

508    recommendations. *Journal of Neurology, Neurosurgery and Psychiatry*, **91**, 638–649.

509    https://doi.org/10.1136/jnnp-2019-322180.

510    PRINSEN, C.A.C., MOKKINK, L.B., BOUTER, L.M., ALONSO, J., PATRICK, D.L., DE VET, H.C.W.,

511    and TERWEE, C.B., 2018, COSMIN guideline for systematic reviews of patient-reported

512    outcome measures. *Quality of Life Research*, **27**, 1147–1157.

513    https://doi.org/10.1007/s11136-018-1798-3.

514    RONEN, G.M., ROSENBAUM, P., and STREINER, D.L., 2000, Outcome measures in pediatric

515    neurology: Why do we need them? *Journal of Child Neurology*, **15**, 775–780.

516    https://doi.org/10.1177/088307380001501201.

517    ROULSTONE, S., COAD, J., AYRE, A., HAMBLY, H., and LINDSAY, G., 2013, *The preferred*

518    *outcomes of children with speech, language and communication needs and their parents.*

519    (UK Department for Education Research report).

520    ROYAL COLLEGE OF SPEECH & LANGUAGE THERAPISTS, 2020, Outcome measurement [online].

521    Available: https://www.rcslt.org/speech-and-language-therapy/guidance-for-delivering-slt-

522    services/outcome-measurement [Accessed 11 January 2020].

523    SCHRIBERG, L.D., AUSTIN, D., LEWIS, B.A., MCSWEENEY, J.L., and WILSON, D.L., 1997, The

524    speech disorders classification system (SDCS): extensions and lifespan reference data.

525    *JSLHR*, **40**, 723–40.

526    SPARROW, S., CICCHETTI, D., and BALLA, D., 2005, *Vineland Adaptive Behavior Scales: Survey*

527     *Interview, 2nd edn.* (Minneapolis: Pearson).

528     SPEECH-LANGUAGE & AUDIOLOGY CANADA, 2010, *Position statement on outcomes measures*

529     (Canada: Ottawa).

530     SPEECH-LANGUAGE & AUDIOLOGY CANADA, 2012, SAC Position Paper: Early Identification of

531     Speech & Language Disorders [online]. Available: http://www.sac-oac.ca/professional-

532     resources/resource-library/early-identification-speech-language-

533     disorders?_ga=1.225121811.849982984.1468963630.

534     SQUIRES, J., BRICKER, D., and TWOMBLY, E., 2003, *The ASQ-SE User's Guide: Ages and Stages*

535     *Questionnaires Social-Emotional. A Parent-Completed, Child-Monitoring Systemfor Social-*

536     *Emotional Behaviours.* (Baltimore: Paul H. Brookes Publishing Co.).

537     TATE, R.L., 2019, Measuring outcomes and monitoring progress in the era of evidence-based

538     clinical practice. *Brain Impairment*, 276–288. https://doi.org/10.1017/BrImp.2019.28.

539     TERWEE, C.B., PRINSEN, C.A.C., CHIAROTTO, A., WESTERMAN, M.J., PATRICK, D.L., ALONSO, J.,

540     BOUTER, L.M., DE VET, H.C.W., and MOKKINK, L.B., 2018, COSMIN methodology for

541     evaluating the content validity of patient-reported outcome measures: A Delphi study.

542     *Quality of Life Research*, **27**, 1159–1170. https://doi.org/10.1007/s11136-018-1829-0.

543     THOMAS-STONELL, N.L., CUNNINGHAM, B.J., ROBERTSON, B., and ROSENBAUM, P., 2020, *The*

544     *Focus on the Outcomes of Communication Under Six (FOCUS©) Manual, Second Edition*

545     (Hamilton, ON: CanChild).

546     THOMAS-STONELL, N.L., McCONNEY-ELLIS, S., ODDSON, B., ROBERTSON, B., and ROSENBAUM,

547     P., 2007, An evaluation of the responsiveness of the Pre-Kindergarden ASHA NOMS.

548     *Canadian Journal of Speech-Language Pathology and Audiology*, **31**, 74–82.

549     THOMAS-STONELL, N.L., ODDSON, B., ROBERTSON, B., and ROSENBAUM, P., 2009, Predicted and

550     observed outcomes in preschool children following speech and language treatment: Parent

551    and clinician perspectives. *Journal of Communication Disorders*, **42**, 29–42.

552    https://doi.org/10.1016/j.jcomdis.2008.08.002.

553    THOMAS-STONELL, N.L., ODDSON, B., ROBERTSON, B., and ROSENBAUM, P.L., 2010,

554    Development of the FOCUS (Focus on the Outcomes of Communication Under Six), a

555    communication outcome measure for preschool children. *Developmental Medicine and*

556    *Child Neurology*, **52**, 47–53. https://doi.org/10.1111/j.1469-8749.2009.03410.x.

557    THOMAS-STONELL, N.L., ODDSON, B., ROBERTSON, B., and ROSENBAUM, P.L., 2013, Validation

558    of the Focus on the Outcomes of Communication under Six outcome measure.

559    *Developmental Medicine and Child Neurology*, **55**, 546–552.

560    https://doi.org/10.1111/dmcn.12123.

561    THOMAS-STONELL, N.L., ROBERTSON, B., WALKER, J., ODDSON, B., WASHINGTON, K.N., and

562    ROSENBAUM, P., 2015, *FOCUS©: Focus on the Outcomes of Communication Under Six*

563    *Manual* (Toronto: ON, Canada: Holland Bloorview Kids Rehabilitation Hospital).

564    THREATS, T.T., 2013, WHO's International Classification of Functioning, Disability, and Health:

565    A framework for clinical and research outcomes. In *Outcomes in speech-language*

566    *pathology* ((New York, NY)), pp. 58–72.

567    VAN TIGGELEN, H., KOTTNER, J., CAMPBELL, K., LEBLANC, K., WOO, K., VERHAEGHE, S., VAN

568    HECKE, A., and BEECKMAN, D., 2020, Measurement properties of classifications for skin

569    tears: A systematic review. *International Journal of Nursing Studies*, **110**, 103694.

570    https://doi.org/10.1016/j.ijnurstu.2020.103694.

571    VALLINO-NAPOLI, L.D. and REILLY, S., 2004, Evidence-based health care: A survey of speech

572    pathology practice. *International Journal of Speech-Language Pathology*, **6**, 107–112.

573    https://doi.org/10.1080/14417040410001708530.

574    VARNI, J.W., 1998, The PedsQL Measurement Model for the Pediatric Quality of Life Inventory

575    [online]. Available: http://www.pedsql.org/%0Aabout_pedsql.html [Accessed 3 March

576    2020].

577    WASHINGTON, K.N., ODDSON, B., ROBERTSON, B., ROSENBAUM, P., and THOMAS-STONELL, N.,

578    2013, Reliability of the Focus on the Outcomes of Communication Under Six (FOCUS©).

579    *Journal of Clinical Practice in Speech-Language Pathology*, **15**, 25–31.

580    https://doi.org/10.1111/cch.12049.

581    WASHINGTON, K.N., THOMAS-STONELL, N., ODDSON, B., MCLEOD, S., WARR-LEEPER, G.,

582    ROBERTSON, B., and ROSENBAUM, P., 2013, Construct validity of the FOCUS© (Focus on

583    the Outcomes of Communication Under Six): A communicative participation outcome

584    measure for preschool children. *Child: Care, Health and Development*, **39**, 481–489.

585    https://doi.org/10.1111/cch.12043.

586    WILCOX, K. and MORRIS, S., 1999, *Children's Speech Intelligibility Measure. 1999* (San Antonio:

587    The Psychological Corporation Harcourt Brace & Company,).

588    WILLIAMS, B. and BEOVICH, B., 2020, A systematic review of psychometric assessment of the

589    Jefferson Scale of Empathy using the COSMIN Risk of Bias checklist. *Journal of*

590    *Evaluation in Clinical Practice*, **26**, 1302–1315. https://doi.org/10.1111/jep.13293.

591    WORLD HEALTH ORGANIZATION, 2001, *International classification of functioning, disability and*

592    *health: ICF*. (Geneva).

593