

# Accepted Manuscript

An emotional functioning item bank of 24 items for computer adaptive testing (CAT) was established

Morten Aa. Petersen, Eva-Maria Gamper, Anna Costantini, Johannes M. Giesinger, Bernhard Holzner, Colin Johnson, Monika Sztankay, Teresa Young, Mogens Groenvold, on behalf of the EORTC Quality of Life Group

PII: S0895-4356(15)00421-7

DOI: [10.1016/j.jclinepi.2015.09.002](https://doi.org/10.1016/j.jclinepi.2015.09.002)

Reference: JCE 8975

To appear in: *Journal of Clinical Epidemiology*

Received Date: 15 December 2014

Revised Date: 4 September 2015

Accepted Date: 7 September 2015

Please cite this article as: Petersen MA, Gamper E-M, Costantini A, Giesinger JM, Bernhard Holzner Colin Johnson Sztankay M, Young T, Groenvold M, on behalf of the EORTC Quality of Life Group, An emotional functioning item bank of 24 items for computer adaptive testing (CAT) was established, *Journal of Clinical Epidemiology* (2015), doi: 10.1016/j.jclinepi.2015.09.002.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



1 **An emotional functioning item bank of 24 items for computer adaptive testing**  
2 **(CAT) was established**

3

4 **Morten Aa. Petersen<sup>a,\*</sup>, Eva-Maria Gamper<sup>b</sup>, Anna Costantini<sup>c</sup>, Johannes M. Giesinger<sup>b</sup>,**  
5 **Bernhard Holzner<sup>b</sup>, Colin Johnson<sup>d</sup>, Monika Sztankay<sup>b</sup>, Teresa Young<sup>e</sup> & Mogens**  
6 **Groenvold<sup>a,f</sup> on behalf of the EORTC Quality of Life Group**

7

8 <sup>a</sup> The Research Unit, Department of Palliative Medicine, Bispebjerg Hospital, University of  
9 Copenhagen , Copenhagen, Denmark

10 <sup>b</sup> Department of Psychiatry and Psychotherapy, Innsbruck Medical University, Innsbruck, Austria

11 <sup>c</sup> Psychoncology Unit, Sant'Andrea Hospital, Faculty of Medicine and Psychology Sapienza  
12 University, Rome, Italy

13 <sup>d</sup> Surgical Unit, University of Southampton, Southampton, UK

14 <sup>e</sup> Lynda Jackson Macmillan Centre, Mount Vernon Cancer Centre, Northwood, Middx, UK

15 <sup>f</sup> Institute of Public Health, University of Copenhagen, Copenhagen, Denmark

16

17 \* Corresponding author: The Research Unit, Department of Palliative Medicine, Bispebjerg  
18 Hospital, Bispebjerg bakke 23, 2400 Copenhagen NV, Denmark. Telephone: (+45) 3531 2025. Fax:  
19 (+45) 3531 2071. Email: [Morten.Aagaard.Petersen@regionh.dk](mailto:Morten.Aagaard.Petersen@regionh.dk)

20

21

22

23

**Abstract**

**Objective:** To improve measurement precision the EORTC Quality of Life Group is developing an item bank for computerized adaptive testing (CAT) of emotional functioning (EF). The item bank will be within the conceptual framework of the widely used EORTC QLQ-C30 questionnaire.

**Study Design and Setting:** Based on literature search and evaluations by international samples of experts and cancer patients 38 candidate items were developed. The psychometric properties of the items were evaluated in a large international sample of cancer patients. This included evaluations of dimensionality, IRT model fit, differential item functioning (DIF), and of measurement precision/statistical power.

**Results:** Responses were obtained from 1,023 cancer patients from four countries. The evaluations showed that 24 items could be included in a unidimensional IRT model. DIF did not seem to have any significant impact on the estimation of EF. Evaluations indicated that the CAT measure may reduce sample size requirements by up to 50% compared to the QLQ-C30 EF scale without reducing power.

**Conclusion:** Based on thorough psychometric evaluations we have established an EF item bank of 24 items. This will allow for more precise and flexible measurement of EF, while maintaining backward compatibility with the QLQ-C30 EF scale.

**Key words:** Computer adaptive test; EORTC QLQ-C30; emotional functioning; item response theory; oncology; patient-reported outcome

**Running head:** An item bank of 24 items for CAT measurement of emotional functioning

**Word count:** 4873

## 49 1. Introduction

50 Computerized adaptive testing (CAT) is a form of intelligent questionnaire; the basic idea is to  
51 maximize the precision by only asking questions relevant for the individual [1-3]. For example, if a  
52 patient has reported severe emotional problems to the previous items (questions), the next item will  
53 be one relevant for patients with severe problems. In this sense the questionnaire is adapted “on-the-  
54 fly” to the individual using previous responses to select the most informative next item. Clearly,  
55 such adaptation cannot be done using usual paper questionnaires, but requires the use of computer  
56 technology. All items used in a CAT are selected from a collection of items called an item bank or  
57 item pool. In a CAT item bank the items have been calibrated (fitted) to an item response theory  
58 (IRT) model [4, 5]. This means that scores based on any subset of the items are comparable. This  
59 unique property facilitates the adaptation to the individual without compromising comparability  
60 across individuals. The adaptability, i.e. the selection of the most informative item at each step,  
61 generally makes CAT instruments more precise than traditional, “static” questionnaires asking the  
62 same number of items and more efficient in the sense that fewer items are needed to obtain a  
63 specific precision. CAT instruments are also highly flexible as they can be adapted to the  
64 requirements of each study or setting. Because of these advantages of CAT several groups have  
65 developed and/or explored the use of CAT to assess patient-reported outcomes (PROs) [6-13].

66  
67 The European Organisation for Research and Treatment of Cancer (EORTC) Quality of Life  
68 questionnaire (QLQ-C30) is an internationally widely used instrument for the assessment of health  
69 related quality of life (HRQOL) in cancer patients [14, 15]. It consists of 30 items measuring 15  
70 aspects of HRQOL: five functional measures, nine symptom measures and one measure of overall  
71 health/quality of life [16]. To improve the assessment of PROs in oncology, the EORTC Quality of  
72 Life Group is currently developing CAT versions of the EORTC QLQ-C30 scales [17-23]. The new  
73 CAT instrument operates within the same conceptual framework as the QLQ-C30. Hence, the aim  
74 is to develop a unidimensional item bank for each QLQ-C30 scale, which, in addition to the original  
75 QLQ-C30 items, consists of items covering the same aspects of the dimension as the QLQ-C30.  
76 That is, each new item bank will include all QLQ-C30 items from the relevant scale. To further  
77 enhance the compatibility with the QLQ-C30 and to ensure a homogeneous and simple format, the  
78 new items should have the same item style as the QLQ-C30 items, i.e. they should employ the same  
79 response options and recall period. In this way the CAT instrument will measure some well-  
80 validated and (to many) well-known HRQOL dimensions and it can be related to the substantial  
81 literature of studies using the QLQ-C30.

82

83 One of the key domains of the QLQ-C30 is emotional functioning (EF). The QLQ-C30 EF scale  
84 consists of four items measuring depression, anxiety (two items), and general distress that are  
85 assumed to represent a unidimensional construct. Responses to the four items are summed to form a  
86 unidimensional EF score. The new item bank should include the four QLQ-C30 EF items. Hence,  
87 the aim is a unidimensional item bank comprising the QLQ-C30 EF items and as many additional  
88 items on depression, anxiety, and general distress as possible.

89

90 As for any EORTC instrument, development of the CAT instrument takes place in an international,  
91 cross-cultural setting. The EORTC CAT development procedure consists of four phases: I)  
92 literature search, II) operationalisation, III) pre-testing, and IV) field testing. Phases I to III of the  
93 EF CAT development have been completed and described elsewhere [17]. In phase I we identified  
94 1,729 EF items from existing questionnaires. The large majority of these items (1,480) were  
95 excluded mainly due to redundancy or lack of relevance for the EORTC measurement of EF. The  
96 remaining items formed the basis for formulating new EF items fitting the “QLQ-C30 item style”.  
97 After a second round of evaluations of redundancy and relevance 63 items were retained.  
98 Evaluations by international samples of experts (phase II) and cancer patients (phase III) further  
99 reduced this to 38 candidate items. The present paper reports on the phase IV field testing and  
100 psychometric evaluations of the 38 EF items.

101

## 102 **2. Methods**

103 The methods and analyses used in phase IV for the final development of the EF item bank are  
104 described below. They generally follow the approach previously reported for other dimensions [20-  
105 22]. Please refer to these publications for further details.

106

### 107 *2.1. Sample*

108 The EORTC CAT is intended for international use for cancer patients in general. Therefore, we  
109 accrued an international sample of cancer patients with different diagnoses, stages of disease, etc.  
110 Patients were recruited from oncology departments in Austria, Denmark, Italy, and the UK in the  
111 period February to December 2011. Patients were invited either by mail or when coming to the  
112 department. Eligibility requirements included a verified cancer diagnosis, age at least 18 years, and  
113 being physically and mentally competent to complete the questionnaire. Written informed consent  
114 was obtained following local requirements.

115

116 The study was approved by the local ethics committees of the participating countries.

117

## 118 2.2. Questionnaire

119 Similar to the QLQ-C30 EF items, all developed items have the recall period “during the past week”  
120 and employ a 4-point response scale: “not at all”, “a little”, “quite a bit” and “very much”. The  
121 QLQ-C30 EF items are about distress and are formulated in the form “did you have this problem”.  
122 The majority of the new items were formulated similarly covering the emotional aspects depression,  
123 anxiety and general distress/negative affect [17]. However, to try to also capture positive emotional  
124 states (i.e. emotional wellbeing/positive affect) and thereby possibly extend the measurement range,  
125 five items were formulated positively, e.g. “Have you felt cheerful?”. The QLQ-C30 EF scale is  
126 scored so that higher values reflect better emotional functioning (less distress). To be in line with  
127 this scoring, we reversed the distress items for the analysis, so that higher scores reflected better  
128 functioning for all items. In addition to the 38 EF items we collected information to clarify whether  
129 the patients found any of the questions problematic and information on the patients’  
130 sociodemographic status.

131

## 132 2.3. Analysis plan

133 The psychometric evaluations and the selection of items were organized into seven steps:

134

135 *1. Descriptive and basic statistical analyses:* This included response rates, item means and standard  
136 deviations, and correlations with the original 4-item QLQ-C30 EF sum scale.

137

138 *2. Assessing patient feedback on the items:* The patients’ qualitative comments and their ratings of  
139 the items were used to assess whether some of the items seemed problematic (e.g. difficult to  
140 understand or confusing).

141

142 *3. Evaluation of dimensionality and local dependence:* As the QLQ-C30 EF scale is  
143 unidimensional, the new CAT based EF measure is also intended to be unidimensional. We  
144 investigated the dimensionality of the items using factor analysis methods for ordinal categorical  
145 data [24]. This included evaluations of dimensionality using scree plot [25], parallel analysis [26,  
146 27] and the Hull method [28]. These analyses were performed using the specialized software  
147 FACTOR v. 9.3.1 [29]. Further, we evaluated the fit of a unidimensional model. The following  
148 criteria were used as indication of reasonable model fit: the root mean square error of  
149 approximation (RMSEA) $<0.10$ , the Tucker-Lewis Index (TLI) $>0.90$  and the Comparative Fit Index  
150 (CFI) $>0.90$  [30, 31]. The analyses were performed using Mplus [24].

151

152 Standard IRT models assume that the items are locally independent, i.e. item responses are

153 independent when controlling for the overall level of EF. This was examined using residual  
154 correlations from the factor model. Residual correlations  $<0.20$  were regarded as indication of local  
155 independence [6].

156

157 *4. Calibration of IRT model and evaluation of item fit:* IRT models are a class of statistical models  
158 used to model latent variables, i.e. variables that are not directly observed but rather inferred from  
159 other (observed) variables [5, 18, 19]. We used the generalized partial credit model (GPCM) [32], a  
160 two-parameter IRT model, to form the basis for the EF CAT. In the GPCM each item has a  
161 discrimination/slope parameter describing an item's ability to discriminate between people with  
162 different scores, and a set of threshold parameters, which defines where on the EF continuum  
163 adjacent response options are equally likely to be endorsed. An item is generally most informative  
164 in the vicinity of the thresholds. The average of an item's thresholds is called the item location.

165

166 Standard IRT models assume monotonicity, i.e. the better EF the more likely it should be to give a  
167 response to the item reflecting good EF. This was examined by inspecting the average item scores  
168 across the rest scores (the sum score of all items except the item in question). If an item complies  
169 with monotonicity, the average item score should not decrease for increasing values of the rest score  
170 [33].

171

172 We used Parscale for estimating the IRT model [34]. Item fit was examined using the item-fit test  
173  $S-X^2$  proposed by Orlando and Thissen [35] and implemented for polytomous items in the SAS  
174 macro IRTFIT [36]. The performance of  $S-X^2$  has been found to be superior to other fit indices [35,  
175 37, 38]. Furthermore, we calculated the average difference between expected and observed item  
176 responses (bias) and the infit and outfit indices [39]. Infit and outfit are both based on mean square  
177 residuals, but assess slightly different aspects of fit. The infit gives more weight to responses from  
178 respondents with an EF score close to the item's location, whereas the outfit is more sensitive to  
179 unexpected responses from respondents far from the item's location. Here the infit is particularly  
180 important, since respondents in CAT measurement are primarily asked items with a location close  
181 to their EF score. For both indices values between 0.7 and 1.3 are often regarded as acceptable [40].  
182 Large values ( $>1.3$ ) indicate misfit to the model while small values ( $<0.7$ ) indicate "overfit", i.e.  
183 better fit than expected statistically, e.g. because of redundancy. Misfit is our main concern, hence,  
184 large values are mainly regarded as problematic.

185

186 *5. Test for differential item functioning (DIF):* DIF analysis evaluates whether the items are  
187 perceived and "behave" similarly in different groups of patients [41]. If this is not the case the item

188 is said to exhibit DIF, which may make comparisons across groups problematic as the same item  
189 response may reflect different levels of EF in the groups. We tested for DIF using ordinal logistic  
190 regression methods [22, 42, 43] with regard to gender, age, country, cancer site, cancer stage,  
191 current treatment, education, work, and cohabitation. Each item was entered as the outcome and the  
192 group (DIF) variables were tested as independent variables controlling for the EF score estimated  
193 using the IRT model calibrated in the previous step.

194  
195 Because of a large sample and multiple testing a difference was regarded as significant if  $p < 0.001$   
196 and potentially relevant if also the coefficient for the group variable (numerically) exceeded 0.64 as  
197 this has been suggested as an indication of moderate to large DIF [43, 44]. For each item, each  
198 group variable was first tested individually for both non-uniform and uniform DIF. To eliminate  
199 false positive DIF findings caused by confounding of the group variables, the group variables  
200 significant in the individual tests were entered together and tested in a multiple logistic regression  
201 model. Only the results of these “multivariate” DIF analyses are reported.

202  
203 We evaluated the practical impact of the significant DIF findings for estimation of EF using a  
204 method proposed by Hart et al [45] which we have further developed [10, 11]. This method  
205 compares the EF scores obtained with the model from step 4 (which does not account for DIF) with  
206 the scores obtained with an IRT model accounting for DIF. If the EF scores obtained with the two  
207 models differed substantially this was regarded as indicating practically problematic DIF, also  
208 termed “salient scale-level differential functioning” [20, 22, 45]. More precisely, if the EF scores  
209 obtained with the two models differed more than the median standard error for the EF estimates this  
210 was regarded as “salient scale-level differential functioning”.

211  
212 *6. Evaluation of discarded items:* We added the discarded items one at a time to the list of items  
213 obtained after step 5 and evaluated whether the item still showed misfit. If this was not the case, i.e.  
214 if it had erroneously been discarded, it was included again.

215  
216 For the final selection of items we calculated the information function. The information function is  
217 a measure of the measurement precision of an item at different levels of EF. The information of the  
218 individual items can be combined to obtain the total information of the item bank.

219  
220 *7. Evaluation of measurement properties:* The measurement properties of the resulting EF CAT  
221 were evaluated using simulations of CAT administration based on the collected data. We simulated  
222 CATs asking 1, 2, ... up to all but 1 items, respectively, estimated the EF score based on these

223 CATs, and compared these scores with the score based on all items. We evaluated the relative  
224 validity (RV) of these CATs as compared to the QLQ-C30 EF scale in detecting expected group  
225 differences [46]. The RV is the ratio of two test statistics for comparing two (known) groups. We  
226 used the t-test statistic for each of the CATs as the numerator and the t-test for the QLQ-C30 EF  
227 scale as the denominator. An  $RV > 1$  indicates that the CAT measure has greater discriminating  
228 power than the QLQ-C30 scale. We hypothesized that patients with stage I or II vs. stage III or IV,  
229 patients not in treatment vs. patients in treatment, patients working vs. patients not working, men  
230 and older patients would have better emotional functioning. The known groups variables being  
231 significant for at least one of the outcomes (the QLQ-C30 EF scale or one of the CAT based scores)  
232 were used for calculating RVs. We also evaluated the RV of the CATs based on simulated data. We  
233 simulated responses to the items based on EF scores sampled from normal distributions with  
234 different means. We compared groups of size  $N_1=N_2=25, 50, \text{ and } 100$ , respectively and true effect  
235 sizes (ESs) of 0.2, 0.5, and 0.8, respectively. For each of these  $3 \times 3 = 9$  possible settings, we ran  
236 2,000 simulations. For further details on these simulations please see Petersen et al. [21].

237

238 The descriptive analyses in step 1 were based on all available data for each item. The analyses in  
239 step 3-6 were based on complete cases, i.e. those responding to all items, while the observed data  
240 evaluations in step 7 were based on those responding to all items in the final model.

241

>> Insert Table 1. Sociodemographic and clinical characteristics of the study sample (N=1,023).<<

ACCEPTED MANUSCRIPT

243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277

### 3. Results

We obtained responses from 1,023 cancer patients coming from Austria, Denmark, Italy, and the UK. Patient characteristics are presented in >> Insert Table 1.

*1. Descriptive and basic statistical analyses.* Response rates to the 38 EF items were generally high (98.2%-99.5%). The average item scores ranged 1.2-2.8 on a 0-3 scale, with 3="not at all", generally reflecting relatively good emotional functioning/little distress. Polychoric correlations between the 34 new items and the sum scale of the original four QLQ-C30 EF items ranged 0.52-0.81 for the (reversed) distress items while the five (unreversed) positive items correlated 0.32-0.55. The lowest were for "Have you felt that you have inner strengths and abilities?" (0.32) and "Have you achieved satisfaction from things that you did?" (0.33).

*2. Assessing patient feedback on the items.* Each item had been rated problematic (difficult to understand, annoying etc.) by no more than four patients (0.4%) and none of the items had been rated intrusive. Hence, generally the patients did not find the items problematic.

*3. Evaluation of dimensionality and local dependence.* Inspection of eigenvalues from an exploratory factor analysis indicated that 57% of the total variation was explained by the first factor. Three additional factors had eigenvalues above 1, and the second factor explained > 5% of the variation. The scree plot suggested one clearly important factor and three potentially important factors. The parallel analysis indicated two possibly three important factors while the Hull method indicated that one factor might be sufficient to explain the variation in the data (details omitted). In a 3-factor and a 4-factor solution the five positive items had their primary loadings on one factor while all the distress items had their primary loadings on other factors. Although the conclusion varied with the method, taken together, these analyses indicated that it might not be sensible to include all 38 items in a unidimensional model, and that the positive items may measure something distinct from the other items. This was confirmed by poor fit indices for a 1-factor model including all 38 items: RMSEA=0.144, CFI=0.690 and TLI=0.954. Excluding the five positive items improved fit, but it was still not acceptable: RMSEA=0.124, CFI=0.788 and TLI=0.971. Next, among the 33 distress items we deleted the poorest fitting items, one at a time, until a 1-factor model had acceptable fit. This resulted in a 24-item model with RMSEA=0.089, CFI=0.906 and TLI=0.987. One factor explained 65% of the variation for these 24 items. Scree plot, parallel analysis and the Hull method all indicated that one factor was sufficient to explain the variation in the 24 items.

278

279 All 276 residual correlations for the 24 items were  $< 0.20$ , and except one all were  $< 0.15$ . Hence,  
280 these correlations did not indicate any local dependence among the retained items.

ACCEPTED MANUSCRIPT

281

282 >> Insert Table 2. Parameter estimates and fit statistics for the 24 items in the final IRT model. <<

283

ACCEPTED MANUSCRIPT

284

285 *4. Calibration of the IRT model and evaluation of item fit.* There were no indications of violations  
286 of the assumption of monotonicity for the 24 items (details omitted). Therefore, we calibrated a  
287 GPCM to the 24 items and evaluated the item fit (results are summarized in >> Insert Table 2). The  
288 tests of item fit indicated good fit to the model for all items. Estimates of bias (average difference  
289 between expected and observed item responses) were all  $\leq 0.01$ . The infit statistics were between  
290 0.93 and 1.07 indicating good fit, and the outfit statistics ranged 0.59-0.97. Items 14, 20 and 26 had  
291 outfits just below 0.7 indicating that the information from these items might be slightly redundant to  
292 the information from other items. However, as all residual correlations for the three items were  
293  $< 0.15$  the assumption of local independence did not seem significantly violated and as the items  
294 may contain unique information making them preferable in specific situations we retained them in  
295 the item bank. Hence, all in all, the 24 items seemed to have acceptable fit to the GPCM.

296

297 *5. Test for DIF.* >> Insert Table 3 presents the findings of significant DIF. There was no significant  
298 DIF with regard to treatment, cohabitation, or work. Of the 24 items, 12 showed potential problems  
299 with DIF. Most differences were found between countries; seven items showed country DIF while  
300 1-2 items showed possible DIF with regard to age, gender, stage, cancer site, or education.

301

302 We evaluated the impact of the possible DIF for the estimation of EF. These evaluations generally  
303 indicated that the DIF findings had almost no effect when using all 24 items to estimate the EF  
304 score (details omitted). Even in the extreme case where only the DIF item was used to estimate EF,  
305 the DIF findings did not seem to have any significant impact, i.e. there were no indications of  
306 salient scale-level differential functioning regardless of the number of items used for the estimation.  
307 Therefore, we concluded that the DIF findings likely did not have any practically relevant impact on  
308 EF estimation, and therefore, retained all items in the model.

309

310

311

312

313

314

315

316

317

318 >> Insert Table 3. Results of the DIF analysis. Regression coefficients and p-values for the  
319 significant findings of DIF. <<

320  
321 *6. Evaluation of discarded items.* The evaluations of the 14 items discarded in the previous steps  
322 indicated that adding any of these to the model again would result in significantly poorer model  
323 fit/lack of unidimensionality. Therefore, no items were reinstated and the 24 items with parameter  
324 estimates shown in >> Insert Table 2 constitute the final EF item bank.

325  
326 The 24 items consisted of 15 items mainly covering depression related aspects, five items about  
327 anxiety and four about general distress. >>Insert Fig. 1 shows the total information for the 24 items.  
328 For comparison the figure also shows the information of the four QLQ-C30 EF items and the  
329 information of the four most informative items at each point along the EF continuum. For  
330 illustration, the EF scores obtained if reporting “very much”, “quite a bit”, “a little”, or “not at all”  
331 problem, respectively, to all items are shown. Measurement with the entire item bank provides  
332 reliability  $\geq 95\%$  (information  $\geq 20$ ) from -2.6 to 0.1 (about 3 standard deviation units) and  
333 reliability of at least 90% (information  $\geq 10$ ) from -3.0 to 0.6. Hence, the item bank provides  
334 particularly precise measurement for patients having some level of emotional problems, while for  
335 patients having at most “a little” EF problems the item bank is less precise. The four most  
336 informative items provide reliability  $\geq 90\%$  from -2.3 to -0.2 while the four QLQ-C30 items have  
337 reliability  $< 90\%$  for the entire continuum.

338

339 >>Insert Fig. 1. Test information function for the 24 items in the final model, information of the  
340 four EORTC QLQ-C30 EF items, and of the four most informative items, respectively.<<

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356 >>Insert Fig. 2. Correlations between EF scores based on CATs asking 1, 2,..., 23 items,  
357 respectively, and EF scores based on all 24 items.<<

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372 7. *Evaluation of measurement properties.* For CATs of all lengths the mean and median EF score  
373 was very close to the score based on all 24 items (details omitted). Scores based on three or more  
374 items correlated  $>0.90$  with the scores based on all items (see >>Insert Fig. 2). >>Insert Fig. 3  
375 summarizes the results of the known groups comparisons. Contrary to expectation there was no  
376 significant difference between work statuses. Stage, treatment, gender and age showed significant  
377 differences as hypothesised and were therefore used to calculate mean RVs for the observed data.  
378 Both observed and simulated data indicated low power if using only one item. The observed data  
379 indicated that when 2 or more items are asked in the CAT, sample sizes may be reduced by about  
380 20-50% without loss of power as compared to using the QLQ-C30 EF scale. The reduction  
381 generally increased with the length of the CAT with approximately 20% savings with 2 items, 40%  
382 with 6 items and 50% with 18 or more items. The simulated data on the other hand indicated that  
383 using the CAT may only reduce sample size requirements by up to 15%. Hence, there were  
384 significant differences in the expected reductions in sample size requirements based on the observed  
385 and simulated data.  
386

387 >>Insert Fig. 3. The average relative validity (RV) and relative required sample size using CAT  
388 measurement compared to using the QLQ-C30 EF sum scale based on observed and simulated data,  
389 respectively.<<

390

#### 391 **4. Discussion**

392 The overall aim of the EORTC CAT-project is to develop a new version of the EORTC QLQ-C30  
393 with better and more precise measurement of HRQOL. This will hopefully improve the  
394 identification and treatment of patients' symptoms and problems. In this study we have expanded  
395 the existing QLQ-C30 EF measure to an item bank covering a broader range of the EF continuum.  
396 We obtained an EF item bank of 24 items showing good psychometric properties: factor analysis  
397 indicated acceptable unidimensionality and IRT calibration and evaluations showed good fit to a  
398 GPCM. We found some indications of DIF. However, evaluations indicated that even though there  
399 were statistically significant differences these did not have any practically relevant impact on the  
400 estimation of EF, i.e. EF scores based on the item bank can be compared across patients and studies  
401 regardless of patient characteristics.

402

403 The item bank provides high measurement precision for a wide range of EF. However, for patients  
404 having at most "a little" EF problems the item bank may lack precision. The measurement may be  
405 improved at a later stage by adding new items particularly relevant for these patients, although from  
406 a clinical point of view it may rarely be relevant to determine with high precision whether a patient  
407 has "little" or "very little" emotional problems, since neither is likely to need treatment. Compared  
408 to the original QLQ-C30 EF scale the new EF item bank has markedly higher precision across the  
409 entire continuum, indicating a significant improvement in measurement precision. Also if selecting  
410 the four most informative items at each point across the continuum, there is a significant gain in  
411 precision compared to the four QLQ-C30 items; the four "maximum information" items provide at  
412 least 50% more information than the QLQ-C30 items for about 3 standard deviation units.

413

414 The list of candidate items included five items on emotional wellbeing/positive affect. The  
415 evaluations indicated that they did not fit well into a unidimensional model with the distress items  
416 and were therefore excluded. These positive items were "experimental" as the QLQ-C30 emotional  
417 scale covers distress only. Further, the assessment of distress is the main interest in the context of  
418 identifying impairments and intervention needs. This split into an emotional distress and an  
419 emotional wellbeing component seems in line with findings, e.g., for the General Health  
420 Questionnaire [47]. Hence, although we refer to the item bank as measuring "emotional  
421 functioning" to be in line with the terminology of the QLQ-C30, it may be more appropriate to refer

422 to it as a measure of emotional distress.

423 In addition to the positive items, nine items were deleted because of poor fit to a unidimensional  
424 model. Of these, two items on feeling “anxious” and “restless” may have been somewhat  
425 ambiguous as they, besides anxiety, may reflect feeling eager/excited. An item on “felt life was  
426 meaningless” and one on “felt life isn't worthwhile” were also deleted. Of all items tested, they had  
427 the fewest reported problems. The highly skewed response distributions may have affected their fit.  
428 When new data with more patients having severe emotional problems become available it might be  
429 that these items will show more favorable fit. Hence, they may be candidates for inclusion in a  
430 “version 2” of the item bank. However, based on the current sample it does not seem appropriate to  
431 include the two items. Finally, four items on being “furious”, “angry”, “bad-tempered”, and  
432 “impatient” were deleted. Likely all four, and particularly the first three, relate to anger/aggressive  
433 feelings. This aspect of emotional functioning did not seem to fit well with the other aspects and is  
434 therefore, not covered by our item bank (anger is not part of the original QLQ-C30 EF scale either).  
435 This is an example of the limitation of our unidimensionality requirement: anger may be a relevant  
436 aspect to measure, but seemingly cannot be included in a unidimensional measure with the other  
437 aspects, and is therefore not included here. If we want to measure anger, a separate item bank may  
438 be constructed. Anger has also in other studies been modelled as a separate domain when measuring  
439 emotional distress [48].

440  
441 The item bank is dominated by items on depression (constituting 15 of the 24 items). To get a more  
442 even distribution of the content areas more items on anxiety and general distress could be added.  
443 Note, however, that the main reason for the smaller numbers of anxiety and general distress items is  
444 that it was difficult to formulate many distinct, relevant, and fitting items for these areas. But even  
445 though the item bank is dominated by items on depression, the CAT can be programmed to always  
446 include items from all content areas, thereby ensuring coverage of all areas constituting the EORTC  
447 EF construct. As this will keep the content of the new measure as close as possible to the QLQ-C30  
448 EF scale we propose to use such content balance.

449  
450 The CAT generally provided precise measurement when comparing scores based on CATs of  
451 different lengths to scores based on the entire item bank. The known groups comparisons generally  
452 indicated that by asking two or more items, sample sizes may be reduced without loss of power as  
453 compared to the QLQ-C30 scale. However, the predicted reduction in sample size requirements  
454 using the new measure differed significantly between the observed data and simulated data  
455 analyses; the simulated data indicated a maximum of 15% reduction while the observed data  
456 indicated up to 50% reduction. This suggests that the actual gain may vary across studies depending

457 on patient samples etc. and underline that detailed evaluations of the power of the CAT measure in  
458 independent data are needed. A validation project has been initiated to evaluate the measurement  
459 properties of the EORTC CAT instrument. This validation study will include new countries,  
460 including countries from Eastern Europe and from outside Europe, allowing for evaluation of the  
461 generalisability of the current findings. Although this validation is not completed, the current,  
462 preliminary version the EORTC CAT may be used by other researchers. For more information on  
463 this preliminary use of the EORTC CAT please visit <http://groups.eortc.be/qol/eortc-cat>.

464  
465 The uncertainty of the increased power using the CAT illustrates that although there are clear  
466 theoretical advantages of CAT it is not obvious how these translate into practical gain like increased  
467 power. CAT is more complex to develop, use and understand than traditional questionnaires. When  
468 is this additional complexity worthwhile, and when may a simpler tool suffice? More studies on the  
469 practical advantages of CAT for PRO measurement are relevant.

470  
471 In conclusion, we have developed an item bank of 24 items for CAT measurement of emotional  
472 functioning. The CAT measure was developed in an international setting targeted at cancer patients  
473 in general in a multitude of languages, but it may also be applied to other patients (and the general  
474 population) as well. It will be backward compatible with the QLQ-C30 and hence with the many  
475 studies that have used this questionnaire. The item bank showed good psychometric properties and  
476 high measurement precision, particularly for patients with some degree of emotional problems.  
477 Evaluations indicated that sample sizes may be reduced up to 50% without loss of power, compared  
478 to the QLQ-C30. However, these evaluations were subject to some uncertainty and the  
479 measurement properties should be validated with new data before drawing any final conclusions.

#### 482 **Acknowledgments**

483 The study was funded by grants from the EORTC Quality of Life Group. The work of Eva-Maria  
484 Gamper and Johannes M. Giesinger was funded by a grant from the Austrian Science Fund (FWF  
485 L502 and FWF J3353).

486 The authors would like to thank the patients responding to our items and our collaborators for  
487 collecting these essential patient responses.

488 None of the authors have any conflicts of interest that might have biased the work.

489

490

491

## References

492

493 1. Wainer H. Computerized Adaptive testing: A Primer (2nd). Mahwah, New Jersey: Lawrence  
494 Erlbaum Associates, Inc.; 2000.

495 2. van der Linden WJ, Glas CAW. Computerized Adaptive testing: Theory and Practice. The  
496 Netherlands: Kluwer Academic Publishers; 2000.

497 3. van der Linden WJ, Glas CAW. Elements of Adaptive Testing. New York: Springer; 2010.

498 4. Hambleton RK, Swaminathan H, Rogers HJ. Fundamentals of Item Response Theory. Newbury  
499 Park: Sage Publications, Inc; 1991.

500 5. van der Linden WJ, Hambleton RK. Handbook of Modern Item Response Theory. Berlin:  
501 Springer-Verlag; 1997.

502 6. Bjorner JB, Kosinski M, Ware JE, Jr. Calibration of an item pool for assessing the burden of  
503 headaches: an application of item response theory to the headache impact test (HIT).  
504 Quality of Life Research 2003; 12(8):913-33.

505 7. Cella D, Gershon R, Lai JS, Choi S. The future of outcomes measurement: item banking,  
506 tailored short-forms, and computerized adaptive assessment. Qual Life Res. 2007; 16  
507 Suppl 1:133-41.

508 8. Fliege H, Becker J, Walter OB, Bjorner JB, Klapp BF, Rose M. Development of a computer-  
509 adaptive test for depression (D-CAT). Quality of Life Research 2005; 14(10):2277-91.

510 9. Gibbons RD, Weiss DJ, Pilkonis PA, Frank E, Moore T, Kim JB et al. Development of the  
511 CAT-ANX: A Computerized Adaptive Test for Anxiety. Am.J.Psychiatry 2013;

512 10. Haley SM, Ni P, Fragala-Pinkham MA, Skrinar AM, Corzo D. A computer adaptive  
513 testing approach for assessing physical functioning in children and adolescents.

514 Dev.Med Child Neurol. 2005; 47(2):113-20.

- 515 11. Hart DL, Wang YC, Stratford PW, Mioduski JE. Computerized adaptive test for patients  
516 with foot or ankle impairments produced valid and responsive measures of function.  
517 *Quality of Life Research* 2008; 17(8):1081-91.
- 518 12. Jette AM, Haley SM, Ni P, Olarsch S, Moed R. Creating a computer adaptive test  
519 version of the late-life function and disability instrument. *J Gerontol.A Biol.Sci.Med.Sci.*  
520 2008; 63(11):1246-56.
- 521 13. Riley WT, Pilkonis P, Cella D. Application of the National Institutes of Health Patient-  
522 reported Outcome Measurement Information System (PROMIS) to mental health  
523 research. *J.Ment.Health Policy Econ.* 2011; 14(4):201-8.
- 524 14. Aaronson NK, Ahmedzai S, Bergman B, Bullinger M, Cull A, Duez NJ et al. The  
525 European Organization for Research and Treatment of Cancer QLQ-C30: a quality-of-  
526 life instrument for use in international clinical trials in oncology. *J Natl Cancer Inst*  
527 1993; 85(5):365-76.
- 528 15. Fayers P, Bottomley A. Quality of life research within the EORTC-the EORTC QLQ -  
529 C30. European Organisation for Research and Treatment of Cancer. *European Journal of*  
530 *Cancer* 2002; 38 Suppl 4:S125-S133.
- 531 16. Fayers PM, Aaronson NK, Bjordal K, Groenvold M, Curran D, Bottomley A. The  
532 EORTC QLQ-C30 Scoring Manual (third). Brussels: European Organisation for  
533 Research and Treatment of Cancer; 2001.
- 534 17. Gamper EM, Groenvold M, Petersen MAa, Young T, Costantini A, Aaronson N et al.  
535 The EORTC Emotional Functioning computer adaptive test (CAT): phase I-III of a  
536 cross-cultural item bank development. *Psycho-Oncology* 2014; 23:397-403.
- 537 18. Giesinger JM, Petersen MAa, Groenvold M, Aaronson NK, Arraras JI, Conroy T et al.  
538 Cross-cultural development of an item list for computer-adaptive testing of fatigue in  
539 oncological patients. *Health Qual Life Outcomes* 2011; 9(19)
- 540 19. Petersen MAa, Groenvold M, Aaronson NK, Chie W-C, Conroy T, Costantini A et al.  
541 Development of computerised adaptive testing (CAT) for the EORTC QLQ-C30

- 542 dimensions – General approach and initial results for physical functioning. *European*  
543 *Journal of Cancer* 2010; 46:1352-8.
- 544 20. Petersen MAa, Groenvold M, Aaronson NK, Chie W-C, Conroy T, Costantini A et al.  
545 Development of computerized adaptive testing (CAT) for the EORTC QLQ-C30  
546 physical functioning dimension. *Quality of Life Research* 2011; 20(4):479-90.
- 547 21. Petersen MAa, Aaronson NK, Arraras JI, Chie W-C, Conroy T, Costantini A et al. The  
548 EORTC computer-adaptive tests measuring physical functioning and fatigue exhibited  
549 high levels of measurement precision and efficiency. *Journal of Clinical Epidemiology*  
550 2012; 66(3):330-9.
- 551 22. Petersen MAa, Giesinger JM, Holzner B, Arraras JI, Conroy T, Gamper EM et al.  
552 Psychometric evaluation of the EORTC computerized adaptive test (CAT) fatigue item  
553 pool. *Quality of Life Research* 2013; 22(9):2443-54.
- 554 23. Thamsborg LH, Petersen MAa, Aaronson NK, Chie W-C, Costantini A, Holzner B et al.  
555 Development of a lack of appetite item bank for computer-adaptive testing (CAT).  
556 *Support Care Cancer* 2014; 23(6):1541-8.
- 557 24. Muthen LK, Muthen BO. *Mplus User's Guide* (2nd edition). Los Angeles, CA: Muthen  
558 & Muthen; 2002.
- 559 25. Cattell RB. Scree Test for Number of Factors. *Multivariate Behavioral Research* 1966;  
560 1(2):245-76.
- 561 26. Horn JL. A Rationale and Test for the Number of Factors in Factor-Analysis.  
562 *Psychometrika* 1965; 30(2):179-85.
- 563 27. Timmerman ME, Lorenzo-Seva U. Dimensionality Assessment of Ordered Polytomous  
564 Items With Parallel Analysis. *Psychol Methods* 2011; 16(2):209-20.
- 565 28. Lorenzo-Seva U, Timmerman ME, Kiers HAL. The Hull Method for Selecting the  
566 Number of Common Factors. *Multivariate Behavioral Research* 2011; 46(2):340-64.

- 567 29. Lorenzo-Seva U., Ferrando PJ. FACTOR v. 9.3.1 [computer program]. Available at  
568 <http://psico.fcep.urv.es/utilitats/factor> (Accessed May 2015). 2015;
- 569 30. Browne MW., Cudeck R. Alternative Ways of Assessing Model Fit. *Sociological  
570 Methods & Research* 1992; 21(2):230-58.
- 571 31. Kline RB. Principles and practice of structural equation modeling (2nd). New York: The  
572 Guilford Press; 2005.
- 573 32. Muraki E. A Generalized Partial Credit Model. In: van der Linden WJ, Hambleton RK,  
574 Eds. *Handbook of Modern Item Response Theory* Berlin: Springer; 1997:153-68.
- 575 33. Junker BW., Sijtsma K. Latent and manifest monotonicity in item response models.  
576 *Applied Psychological Measurement* 2000; 24(1):65-81.
- 577 34. Muraki E, Bock RD. PARSCALE - IRT based Test Scoring and Item Analysis for  
578 Graded Open-ended Exercises and Performance Tasks. Chicago: Scientific Software  
579 International, Inc.; 1996.
- 580 35. Orlando M., Thissen D. Likelihood-Based Item-Fit Indices for Dichotomous Item  
581 Response Theory Models. *Applied Psychological Measurement* 2000; 24(1):50-64.
- 582 36. Bjorner JB, Smith KJ, Stone C, Sun X. Software: IRTFIT: A macro for item fit and local  
583 dependence tests under IRT models (obtained at:  
584 [http://outcomes.cancer.gov/areas/measurement/irt\\_model\\_fit.html](http://outcomes.cancer.gov/areas/measurement/irt_model_fit.html)). 2007;
- 585 37. Orlando M., Thissen D. Further investigation of the performance of S-X-2: An item fit  
586 index for use with dichotomous item response theory models. *Applied Psychological  
587 Measurement* 2003; 27(4):289-98.
- 588 38. Kang T., Chen TT. Performance of the Generalized S-X(2) Item Fit Index for  
589 Polytomous IRT Models. *Journal of Educational Measurement* 2008; 45(4):391-406.
- 590 39. Bond TG, Fox CM. *Applying the Rasch model: Fundamental measurement in the human  
591 sciences* (2nd edition). New Jersey: Lawrence Erlbaum Associates, Inc.; 2007.

- 592 40. Wright BD, Linacre JM. Reasonable mean-square fit values. *Rasch Measurement*  
593 *Transactions* 1994; 8(3):370.
- 594 41. Holland PW, Wainer H. *Differential Item Functioning*. Hillsdale, NJ: Lawrence Erlbaum  
595 Associates; 1993.
- 596 42. French AW, Miller TR. Logistic regression and its use in detecting differential item  
597 functioning in polytomous items. *J Educ Meas* 1996; 33(3):315-32.
- 598 43. Petersen MAa, Groenvold M, Bjorner JB, Aaronson NK, Conroy T, Cull A et al. Use of  
599 differential item functioning analysis to assess the equivalence of translations of a  
600 questionnaire. *Quality of Life Research* 2003; 12(4):373-85.
- 601 44. Bjorner JB, Kreiner S, Ware JE, Damsgaard MT, Bech P. Differential item functioning  
602 in the Danish translation of the SF-36. *Journal of Clinical Epidemiology* 1998;  
603 51(11):1189-202.
- 604 45. Hart DL, Deutscher D, Crane PK, Wang YC. Differential item functioning was  
605 negligible in an adaptive test of functional status for patients with knee impairments who  
606 spoke English or Hebrew. *Quality of Life Research* 2009; 18(8):1067-83.
- 607 46. Fayers PM, Machin D. *Quality of Life. The assessment, analysis and Interpretation of*  
608 *patient-reported outcomes* (2nd ed.). Chichester: John Wiley & Sons Ltd; 2007.
- 609 47. Stochl J, Jones PB, Croudace TJ. Mokken scale analysis of mental health and well-being  
610 questionnaire item responses: a non-parametric IRT method in empirical research for  
611 applied health researchers. *BMC Med Res Methodol*. 2012; 12:74.
- 612 48. Pilkonis PA, Choi SW, Reise SP, Stover AM, Riley WT, Cella D. Item banks for  
613 measuring emotional distress from the Patient-Reported Outcomes Measurement  
614 Information System (PROMIS(R)): depression, anxiety, and anger. *Assessment*. 2011;  
615 18(3):263-83.  
616  
617

1 Table 1. Sociodemographic and clinical characteristics of the study sample (N=1,023).

		<b>N/mean</b>
Age (mean years)		62 (range 22-88)
Gender	Male	483 (47%)
	Female	540 (53%)
Country	Austria	204 (20%)
	Denmark	205 (20%)
	Italy	94 (9%)
	UK	520 (51%)
Education	0-10 years	376 (37%)
	11-13 years	258 (25%)
	14-16 years	218 (21%)
	>16 years	158 (15%)
Work	Working	322 (31%)
	Retired	564 (55%)
	Other	125 (12%)
Cohabitation	Living with a partner	759 (74%)
	Living alone	244 (24%)
Cancer stage	I-II	456 (45%)
	III-IV	420 (41%)
Cancer site	Breast	130 (13%)
	Gastrointestinal	199 (20%)
	Gynaecological	97 (10%)
	Head and neck	74 (7%)
	Lung	90 (9%)
	Urogenital	104 (10%)
	Other	235 (23%)
	Current treatment	Chemotherapy
	Other treatment	117 (11%)
	No current treatment	486 (48%)

Table 2. Parameter estimates and fit statistics for the 24 items in the final IRT model.

Item	Thresholds			Location	Item fit p-value	Bias	Infit	Outfit	
	Slope	T1	T2						T3
Item 3: Did you feel tense? (from QLQ-C30)	1.62	-1.79	-1.04	0.10	-0.91	0.837	0.01	0.98	0.89
Item 4: Have you felt helpless?	2.41	-1.67	-1.23	-0.67	-1.19	0.985	0.01	1.00	0.73
Item 5: Have you felt panic?	1.78	-1.80	-1.53	-1.19	-1.51	0.996	0.01	0.99	0.97
Item 6: Have you lost interest in things, such as recreational or social activities (independently of your actual ability to do them)?	1.14	-1.91	-1.21	-0.81	-1.31	0.640	0.01	0.99	0.97
Item 7: Have you felt vulnerable?	2.11	-1.79	-1.24	-0.26	-1.10	0.953	0.01	0.99	0.82
Item 8: Have you felt frustrated?	1.72	-1.65	-1.04	-0.14	-0.94	0.359	0.01	0.99	0.85
Item 9: Have you felt worthless?	2.29	-1.86	-1.48	-1.01	-1.45	0.969	0.01	1.02	0.82
Item 12: Have you felt discouraged?	2.58	-1.93	-1.42	-0.44	-1.27	0.820	0.01	0.97	0.75
Item 13: Have you had emotional outbursts?	1.24	-1.65	-1.60	-0.75	-1.33	0.225	0.01	1.01	0.89
Item 14: Have you felt that nothing could cheer you up?	3.32	-1.95	-1.53	-0.80	-1.43	0.995	0.01	0.98	0.63
Item 15: Have you felt afraid?	1.69	-1.63	-1.53	-0.53	-1.23	0.616	0.01	1.00	0.83
Item 16: Have you felt that pleasure has gone from your life?	1.97	-1.71	-1.47	-0.66	-1.28	0.985	0.01	0.97	0.84
Item 17: Have you had difficulty relaxing?	1.41	-1.78	-1.25	-0.25	-1.10	0.560	0.01	0.99	0.93

Item 18: Have you lost interest in your appearance?	1.33	-2.31	-1.92	-1.33	-1.85	0.543	0.01	1.02	0.93
Item 20: Have you felt miserable?	3.19	-1.75	-1.16	-0.29	-1.06	0.878	0.01	0.93	0.68
Item 22: Did you feel depressed? (from QLQ-C30)	2.53	-1.66	-1.25	-0.29	-1.07	0.946	0.01	0.95	0.82
Item 23: Did you feel irritable? (from QLQ-C30)	1.64	-2.28	-1.50	-0.02	-1.27	0.420	0.01	0.97	0.88
Item 24: Have you felt useless?	2.25	-1.85	-1.45	-0.89	-1.40	0.651	0.01	0.98	0.86
Item 25: Did you worry? (from QLQ-C30)	1.64	-1.66	-0.94	0.44	-0.72	0.718	0.01	0.95	0.92
Item 26: Have you felt desperate?	3.03	-2.10	-1.43	-1.03	-1.52	0.952	0.01	1.00	0.59
Item 29: Have you been afraid of losing control?	1.81	-2.06	-1.71	-1.06	-1.61	0.871	0.01	1.03	0.87
Item 30: Have you felt sad?	2.61	-1.59	-1.14	0.09	-0.88	0.927	0.01	0.93	0.83
Item 31: Have you felt like giving up?	2.17	-1.80	-1.58	-1.37	-1.58	0.992	0.01	1.07	0.71
Item 32: Have you felt that you have nothing to look forward to?*	2.40	-1.52	-0.97		-1.24	0.982	0.01	1.00	0.72

\*: Response options “Quite a bit” and “Very much” were collapsed for item 32 because of reversed thresholds.

Note: The items (text and parameters) constitute the EORTC CAT Emotional Function Item Bank version 1.0. © Copyright 2015 EORTC Quality of Life Group. All rights reserved. User’s Agreements for EORTC Quality of Life Group questionnaires are available from the EORTC Quality of Life Department, Brussels, Belgium, <http://groups.eortc.be/qol/>

Table 3. Results of the DIF analysis. Regression coefficients and p-values for the significant findings of DIF.

Item	DIF	$\beta$	p-value	DIF	$\beta$	p-value
Item 3	Country:	-0.88 (Austria)	<0.0001			
Item 4	Stage:	0.76	<0.0001			
Item 5	No DIF					
Item 6	No DIF					
Item 7	Education:	0.96 (0-10 years)	0.0005			
Item 8	Country:	1.61 (Italy)	<0.0001			
Item 9	No DIF					
Item 12	No DIF					
Item 13	Age:	-1.49 (<40 years)	<0.0001			
Item 14	No DIF					
Item 15	Site:	0.80	0.0001			
Item 16	Country:	-1.61 (Austria)	<0.0001	Gender:	-0.69	<0.0001
Item 17	No DIF					
Item 18	No DIF					
Item 20	Country:	1.16 (Denmark)	<0.0001			
Item 22	Country:	-1.01 (Austria)	<0.0001			
Item 23	No DIF					
Item 24	No DIF					
Item 25	No DIF					

Item 26	Country:	-1.76 (Austria)	<0.0001
Item 29	No DIF		
Item 30	No DIF		
Item 31	Country:	-1.42 (Italy)	<0.0001
Item 32	Gender:	-0.73	0.0001

For country the largest regression coefficient (and the country) when comparing with UK is shown. A coefficient >0 indicates UK patients are more likely to report problems on the item (when controlling for the EF score).

For age the largest regression coefficient (and age group) when comparing with patients  $\geq 70$  years is shown. A coefficient >0 indicates that patients  $\geq 70$  years are more likely to report problems on the item.

For gender a coefficient >0 indicates that women are more likely to report problems on the item.

For cancer site a regression coefficient >0 indicates that breast cancer patients are more likely to report problems.

For cancer stage a regression coefficient >0 indicates that patients with advanced cancer (III-IV) are more likely to report problems.

For education the largest regression coefficient (and group, years of education) when comparing with patients having more than 16 years of education is shown.

A coefficient >0 indicates patients having more than 16 years of education are more likely to report problems on the item.

Fig. 1. Test information function for the 24 items in the final model, information of the four EORTC QLQ-C30 EF items, and of the four most informative items, respectively.

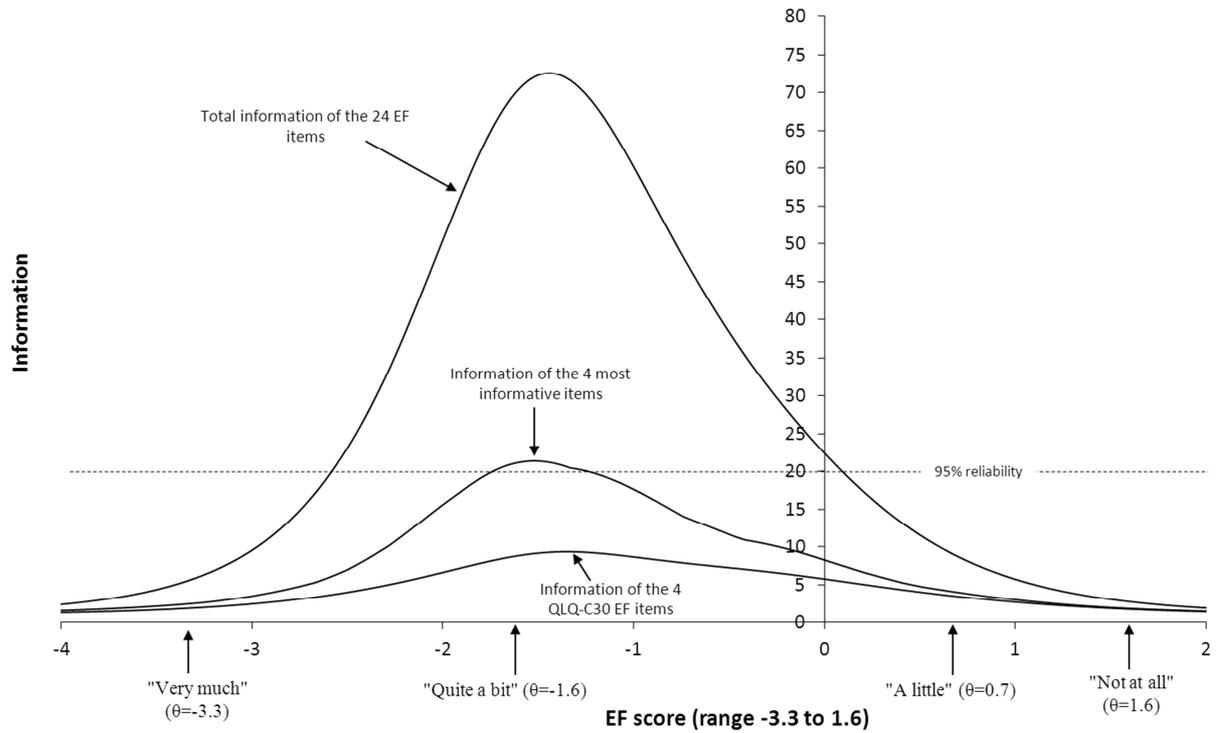


Fig. 2. Correlations between EF scores based on CATs asking 1, 2, ..., 23 items, respectively, and EF scores based on all 24 items.

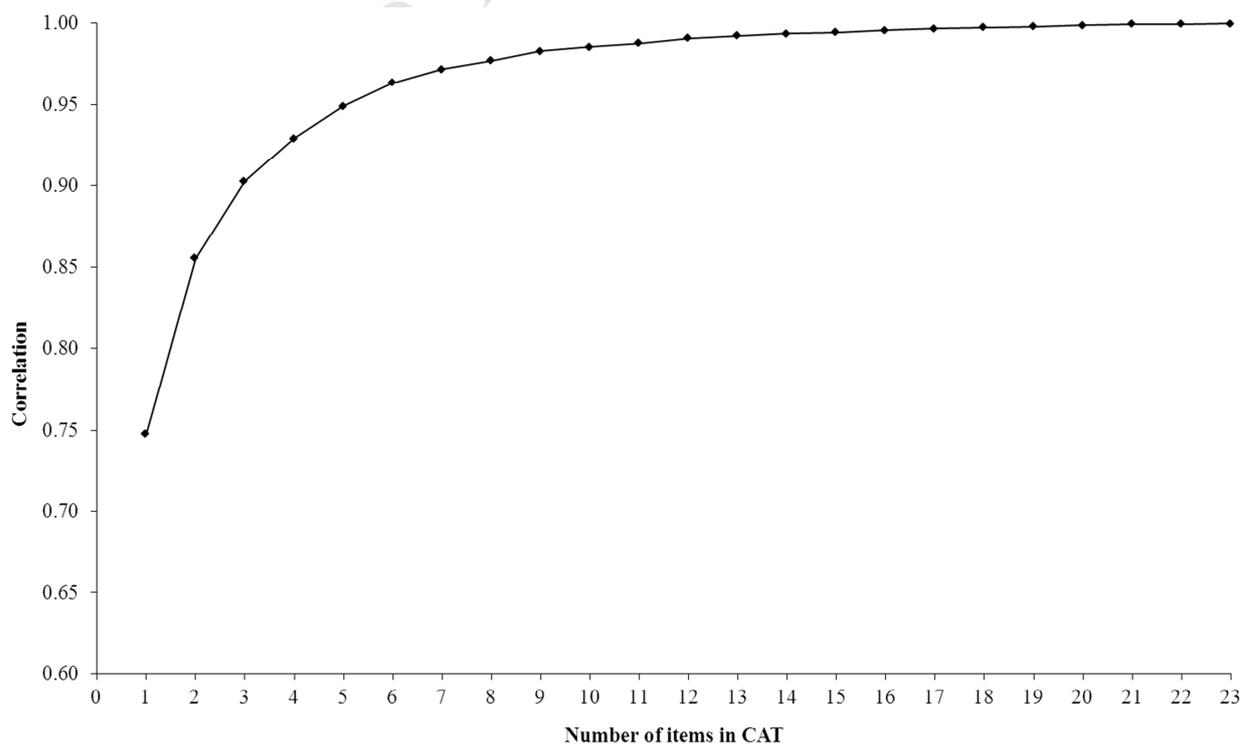
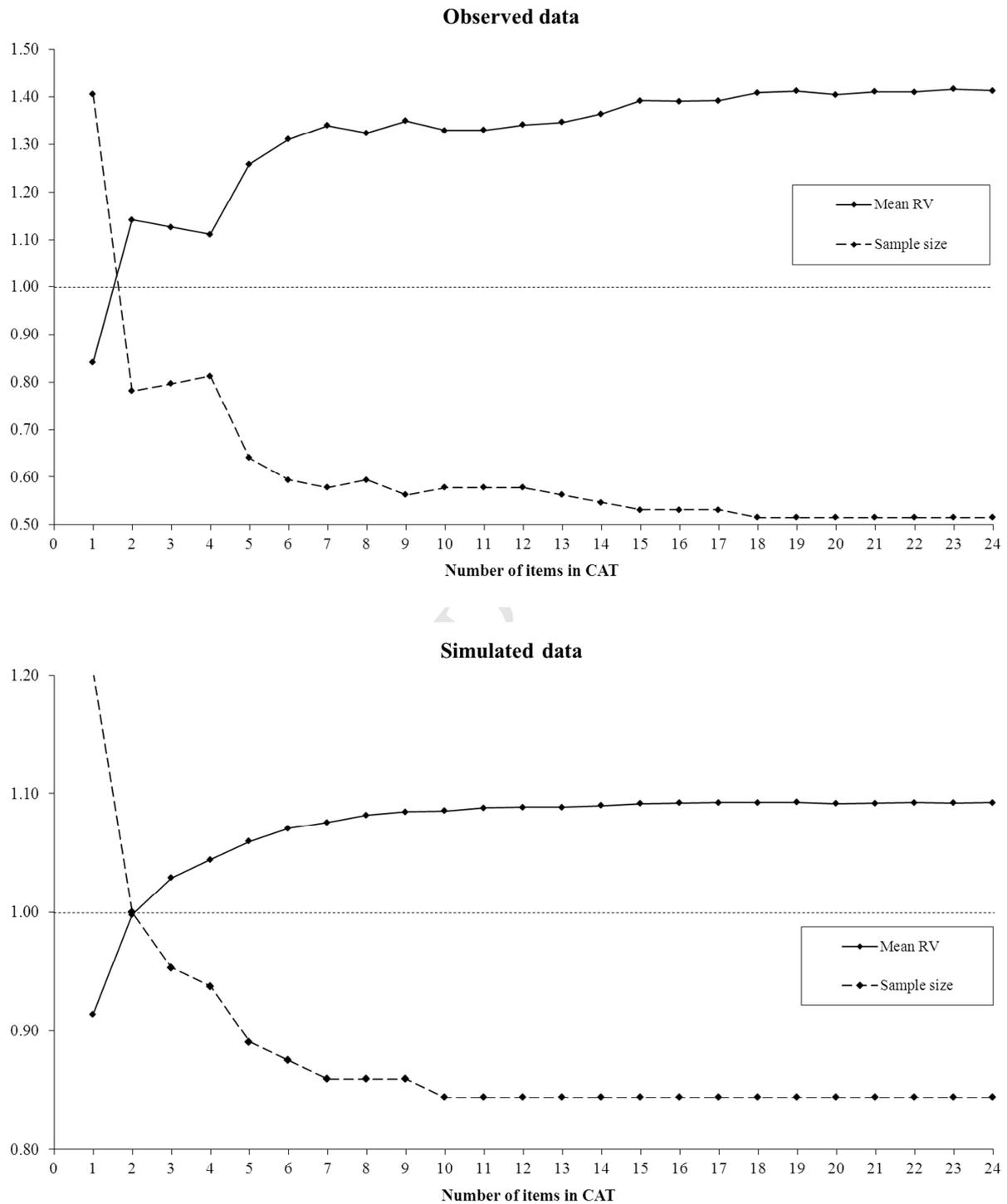


Fig. 3. The average relative validity (RV) and relative required sample size using CAT measurement compared to using the QLQ-C30 EF sum scale based on observed and simulated data, respectively.



Footnote to Fig. 3: For both RV and sample size the plots show the ratio of using the CAT compared to using the QLQ-C30 sum scale. For example, using a CAT with six items the observed data indicates that the validity of the CAT is 1.31 times that of the QLQ-C30 scale ( $RV=1.31$ ) resulting in that only 0.59 (59%) of a sample size used with the QLQ-C30 scale is required when using a 6-item CAT to obtain the same power.

ACCEPTED MANUSCRIPT