1 **Inter-Rater Reliability at the Top End – Measures of Pilots' Non-Technical**

2 **Performance**

3 Running head: Inter-Rater Reliability at the Top End

4 Patrick Gontar & Hans-Juergen Hoermann

5 Abstract

6 **Objective.** The aim of this study is to analyze influences on inter-rater reliability and

7 within-group agreement within a highly experienced rater group when assessing pilots'

8 non-technical skills.

9 **Background.** Non-technical skills of pilots are crucial for the conduct of safe flight

10 operations. To train and assess these skills, reliable expert ratings are required.

11 Literature shows to some degree that inter-rater reliability is influenced by factors

12 related to the targets, scenarios, rating tools, or to the raters themselves.

13 **Method.** Thirty-seven type-rating examiners from a European airline assessed the

14 performance of four flight crews based on video recordings using LOSA and adapted

15 NOTECHS tools. We calculated $r_{wg}$ and *ICC(3)* to measure within-group agreement and

16 inter-rater reliability.

17 **Results.** The findings indicated that within-group agreement and inter-rater reliability

18 were not always acceptable. Both metrics showed that outstanding pilots'

19 performance was rated with higher within-group agreement. For cognitive aspects of

20 performance, inter-rater reliability was higher than for social aspects of performance.

21 Agreement was lower on the pass/fail level than for the distinguished performance

22 scales.

23 **Conclusion.** These results suggest to back pass/fail decisions not exclusively on non-

24 technical skill ratings. We furthermore recommend that regulatory authorities more

25 systematically address inter-rater reliability in airline instructor training. Airlines as

26 well as training facilities should be encouraged to demonstrate sufficient inter-rater

27 reliability when using their rating tools.

28    *Keywords:* inter-rater reliability, within-group agreement, non-technical skills,

29    NOTECHS, LOSA

30

31

## Introduction

33    In-depth accident investigations in the 1970s highlighted the fact that the non-technical

34    behaviors of pilots, like leadership, communication, teamwork and decision making, had clearly been

35    neglected as significant factors for safe flight operations (Cooper, White, & Lauber, 1980). In

36    succession, several approaches attempted to systematically include cockpit/crew resource

37    management (CRM) in pilot training (Helmreich, Merrit, & Wilhelm, 1999). For the evaluation of

38    training success, CRM-skills assessment became relevant. Goldsmith and Johnson (2002) name three

39    major reasons why such an evaluation is important and how it can improve pilot performance: to

40    judge if the *pilot is proficient* enough to fly in the respective airline, to give *sufficient and appropriate*

41    *performance feedback* to the pilot, and to develop and modify the *airline's training program*.

42    Regulatory authorities have provided standards and guidelines for the instruction and assessment of

43    CRM by the airlines (cf. European Aviation Safety Agency, 2011, 2014; European Commission, 2011;

44    Federal Aviation Administration, 2004; Joint Aviation Authorities, 2001). In this context, Robert

45    Helmreich and his colleagues at the University of Texas were very influential, developing behavioral

46    marker systems and other observation methods such as the Line/LOS Checklist for the aviation

47    industry (Helmreich, Klinect, Wilhelm, & Jones, 1999). For European airlines, a rating system called

48    NOTECHS became the standard (Flin et al., 2003; O'Connor, Hoermann, Flin, Lodge, & Goeters, 2002).

49    Because behavioral marker and rating systems are subject to observation bias, aspects of inter-rater

50    reliability (IRR) and inter-rater agreement (IRA) became important (Brannick & Prince, 1991;

51    Brannick, Prince, & Salas, 2002). To ensure that pilots are trained to a required level of competence,

52    reliability of the competence assessment is a vital precondition (cf. Nunnally & Bernstein, 1994).

53    While "…practical elements [of instructor training] should include the development of

54    specific instructor skills, particularly in the area of teaching and assessing threat and error

55    management and CRM" (European Aviation Safety Agency, 2011, FCL.920, p. 282), high IRR and IRA

56    lead to transparent and traceable ratings and can therefore enhance the feedback during the

57  debriefing, and thus the training quality (Gontar & Hoermann, 2014). Unreasonably harsh or

58  unreasonably lenient ratings can not only lead to economic drawbacks, but also to critical safety

59  consequences (Holt, Hansberger, & Boehm-Davis, 2002). For example, raters using overly harsh

60  standards may give rise to needless additional training costs for the airline. An overly harsh rating in

61  an examination flight could jeopardize the pilot's license and have a negative effect on his or her

62  motivation without reasonable cause. The opposite is the case if the raters have pilots passing an

63  examination although they performed below the required minimum. In the latter case, degrading

64  standards will have safety implications.

65      Studies in aviation  (Brannick et al., 2002; Holt et al., 2002; O'Connor et al., 2002; Williams,

66  Holt, & Boehm-Davis, 1997), in air traffic control (Kontogiannis & Malakis, 2013), and in medical

67  domains (Arora et al., 2011; Beard, Marriott, Purdie, & Crossley, 2011; Cooper, Endacott, & Cant,

68  2010; Dedy et al., 2015; Fletcher et al., 2003; Gale et al., 2010; Mitchell et al., 2012; Sevdalis et al.,

69  2008; Yule et al., 2008; Yule et al., 2009) found that professional raters have different views when

70  rating practitioners on their non-technical skills (NTS). The important research tasks in this context

71  are obviously to identify the conditions under which the views of the raters tend to diverge or

72  converge and to apply the outcomes to the improvement of inter-rater reliability. Based on the

73  aforementioned studies, we divided the factors that influence inter-rater reliability into four major

74  themes (comparable to Brannick et al., 2002). These themes are: *target*-related (e.g. target person's

75  level of performance, target person's position in crew), *scenario & task*-related (e.g. taxiing,

76  emergency procedures, cruise flight, approach), *measurement*-related (e.g. rating dimension, scale

77  level, observable markers, anchors), and *rater*-related (e.g. experience, familiarity with rating tools,

78  motivation). It is pointed out that these themes can also be interdependent.

79  ***Target*-related Influences**

80      Regarding *target*-related influences, O'Connor et al. (2002) reported that captains (CPTs)

81  were rated less accurately than first officers (FOs). Mishra, Catchpole, and McCulloch (2009) found

82  slight differences between targets when analyzing IRA for nurses, surgeons, and anesthetists. Yule et

83  al. (2009) found that it is easier to rate targets who perform very well or very poorly than crews

84  whose performance is in the medium range, since extreme behaviors are normally more salient. As

85  average-performing crews represent the majority of cases in reality, it is very important to train

86  inter-rater reliability when rating those (Yule et al., 2009). In addition, these authors pointed out the

3

87    problem that a target's performance may vary on the same dimension (e.g. communication) during

88    the observation period. In this case, it is hard to decide how to weigh the different characteristics

89    and to arrive at a final grade.

90    *Scenario & Task*-**related Influences**

91    The second identified theme, *scenario & task,* is addressed by O'Connor et al. (2002). They

92    were able to show that the content of flight scenarios and tasks influence inter-rater reliability and

93    identified the crucial aspects in their specific scenarios. One major fact they asserted was the

94    difficulty for the rater to "decide how to separate the behaviors and responsibilities of the two

95    pilots" (O'Connor et al., 2002, p. 282). Yule et al. (2008) conducted a study with six different

96    scenarios and found similar results, but attributed them to the special behaviors of the crews, which

97    they stated were easier to rate. Mitchell at al. (2012) explained differences in inter-rater reliability

98    between the scenarios as being affected by the short and variable duration of the scenarios. They

99    furthermore suggested that the semi-scripted scenarios might have influenced inter-rater reliability

100   due to the varying quality of the actors.

101   *Measurement*-**related Influences**

102   Dedy et al. (2015), Mishra et al. (2009), O'Connor et al. (2002), and Yule et al. (2008) found

103   that within-group agreement and inter-rater reliability also depend on the rated dimension.

104   O'Connor et al. (2002) and Yule et al. (2008) reported that interpersonal skills (e.g. communication)

105   were rated in higher agreement than cognitive skills (e.g. decision making). Yule et al. (2008)

106   attributed this effect to the raters, who only had 2.5 hours of training and were not educated in the

107   underlying cognitive models. In contrast, Yule et al. (2009) found an opposite effect: Cognitive skills

108   were rated in higher agreement than social skills. We found this same effect when analyzing pilots'

109   peer and self-rating behaviors (Gontar & Hoermann, 2014). Social aspects such as communication,

110   leadership, and teamwork were rated with lower inter-rater reliability than cognitive aspects, for

111   example work organization, situation awareness, and decision making (Gontar & Hoermann, 2014).

112   We concluded that this effect was due to the scenario, where the successful technical outcome was

113   strongly related to good decision making. Brannick et al. (2002) analyzed the influence of item

114   generality on reliability and found that interjudge agreement was higher for specific behaviors than

115   for a general assessment of CRM in total.

116    *Rater-*related Influences

117        Regarding the influence of the rater, Hamman and Holt (1997) found that factors such as

118    personal interpretation and motivation influence and bias performance ratings (see Flin & Martin,

119    2001). Yule et al. (2008) argued that rating bias also depends on the expertise of the raters in their

120    specific field. They analyzed the average reliability of different rater groups (in this case: general

121    surgeons vs. orthopedic surgeons) and found variance in the agreement, suggesting that "…surgeons'

122    ratings might be more homogeneous when they are rating scenarios based in their own specialty

123    than when rating other specialties" (p. 552). In 2009, Yule et al. showed that prior rating experience

124    can affect rating standards. They compared novice raters with expert raters and found that novices

125    tend to rate more harshly than experts. As possible reasons for low reliability, Weber, Roth, Mavin,

126    and Dekker (2013) suggested that raters might not recognize the same behaviors, or even if they do,

127    they do not evaluate them equally. In a follow-up study, Weber, Mavin, Roth, Henriqson, and Dekker

128    (2014) analyzed the degree to which raters gave different reasons (justifications) for their grading of

129    pilots' behavior. They clustered similar justifications into *topics* and were able to show that raters use

130    different *topics* to assess specific performance categories.

131    **Research Needs**

132        Even though IRR of performance ratings is influenced by the above mentioned factors, check

133    and training practices of the airlines have to rely on instructor pilots assessing the technical and non-

134    technical skills of their trainees. This is primarily done through observation. Several studies have

135    attempted to improve reliability and structural validity of the crew assessment by intensified

136    instructor training (e.g. Holt et al., 2002) or by improving tools for rating non-technical skills (e.g.

137    Sevdalis et al., 2008). Holt et al. (2002) looked at the development of IRR over a period of three years

138    with rater training. They found generally acceptable results, but could not identify strong

139    improvement over the years. These authors mentioned that due to the small number of raters in the

140    beginning, turnover in the group of raters may have affected the group's rating performance.

141    Furthermore, they noted that the rated scenarios differed from year to year. Sevdalis et al. (2008)

142    analyzed the IRR of raters after revising their NTS rating tool. Even after these revisions, specific

143    dimensions such as *cooperation and team skills* showed barely adequate reliability. However, they

144    could not rule out that a lack of familiarity with the revised definitions led to lower IRR for that

145    dimension.

146    In contrast to the studies cited here, this study kept the influence of the *raters* and the

147    influence of the *scenario & task* constant, which allowed us to focus on the influence of *target* and

148    *measurement*. Raters and scenarios were kept unchanged by selecting a homogeneous and very

149    experienced group of type-rating examiners from the same airline and showing them videos with

150    different flight crews performing the same flying tasks. We examined how reliably these raters, who

151    worked for the same airline as the crews, used different rating tools to assess the pilots.

152    To our knowledge, no study has been published in which a large group of raters with

153    homogeneous experience, education, and affiliation participated, in order to keep *rater*-induced

154    effects constant. Furthermore, no study was found that kept the influence of *scenarios & tasks*

155    constant across different crews. In this study we asked the instructor pilots to assess actual flight

156    crews from the same airline in realistic simulator scenarios containing the same task for each crew.

157    Such a situation is very common in reality: All pilots in an airline have to fly the same missions in

158    training and examination flights. In addition, most of the previous studies had either volunteer raters

159    (e.g. Mitchell et al., 2012; Yule et al., 2009) or did not specify how raters had been recruited (e.g.

160    Fletcher et al., 2003; O'Connor et al., 2002). We suggest that using volunteers, and thus self-selected

161    raters, would potentially bias the ratings and therefore would not reflect the daily practice. Normally,

162    instructors and pilots are assigned to their specific training or check missions.

163    In our study, we expect similar results for within-group agreement as reported by O'Connor

164    et al. (2002). These authors found an average $r_{wg}$ of .76 across all the different rating dimensions at

165    the category level of NOTECHS, which is comparable to our $NTS_{dim}$ measurement (see dependent

166    measures). O'Connor et al. (2002) showed that the agreement varied for different scenarios (from

167    .64 to .87). In our study, the scenario remained unchanged. However, different crews exhibited the

168    full range of performance, from *outstanding* to *poor*. We expect that raters show higher agreement

169    for extreme performance than for average performance, because extreme performance is assumed

170    to be more salient (Yule et al., 2009).

171    To summarize the above, the research questions addressed by this paper arise from the two

172    major themes that influence inter-rater reliability: *target* and *measurement*. With respect to the

173    *target*-related influences, we investigate differences in the ratings for the two crew members (CPT

174    and FO) in relation to their level of performance. In terms of the *measurement*-related factor, we

175    analyze the influence of the familiarity with the tools, the tools' dimensions, and the scale levels of

176 the tools. We keep the influences from the *raters* and the *scenario & task* constant by choosing the

177 best raters and using the same flight scenarios for all pilots.

178 **Method**

179 **Participants**

180     A sample of 37 type-rating examiners (TREs) from a major European airline, all holding valid

181 licenses for the Airbus A320, took part as raters in this experiment. Their participation was not

182 voluntary, since they were assigned to this rating experiment as part of a workshop. Due to their

183 specific training and certification, the examiners are the most experienced instructors for this aircraft

184 type within this company. They represent a homogeneous group with regard to their affiliation and

185 experience. The mean age of the participants was 49.9 years ($SD$ = 4.2 years). They had a mean

186 experience of 11.5 ($SD$ = 4.3) years as training and check pilots, and had a mean number of 13,604.2

187 ($SD$ = 3,900.1) airline flight hours.

188     As part of their initial and recurrent instructor courses, all participants had received several

189 days of theoretical and practical training for their rating skills. Rating exercises were done with video

190 examples in classrooms as well as during real training sessions in the simulator. As part of the

191 training, instructors received feedback on their individual rating tendencies. In addition, they

192 participated in annual standardization meetings, which contain specific case study exercises. Since

193 the simulator scenario was new and not previously included in routine recurrent trainings by the

194 airline, none of the raters had specific experience with the presented simulator scenarios – neither as

195 participating pilot nor as instructor pilot.

196 **Apparatus**

197     The 37 raters assessed videotapes of the same flight scenario flown by four different crews in

198 an A320 simulator; the different videotapes were presented on a screen using a projector in a

199 classroom to all raters at the same time. For the purpose of de-identification, the pilots' voices in the

200 videotapes were modified by changing the pitch; dialogs were still clearly understandable.

201 **Recorded Simulator Mission**

202    The presented videotapes were extracted from a mission in a flight simulator study

203    conducted with (*n* = 60) short-haul pilots on the Airbus A320 (see Gontar & Hoermann, 2014; Gontar,

204    Hoermann, Deischl, & Haslbeck, 2014). The flight simulator mission aimed to analyze the pilots'

205    behavior in unforeseen situations under high workload. Therefore, the pilots showed authentic non-

206    scripted behaviors. This experiment was conducted in a full-flight simulator (*JAR STD 1A Level D*), but

207    was not part of pilots' recurrent training within the airline.

208    The selected videotapes for the inter-rater reliability study show a sequence with high task

209    load for the pilots. In the flight simulator mission, the crews began a visual approach (VOR B) to

210    runway 22R at Nice Côte d'Azur Airport (LFMN), 15 miles east of the airport (D15 AZR) at an altitude

211    of 3,000 feet with a speed of 170 knots, and a heading of 269 degrees; there was light rain, the

212    runway was wet, visibility was 10,000 meters, wind was 10 knots from the south, and the

213    temperature was 12° Celsius. The aircraft had 2,500 kg fuel on board (corresponding to a remaining

214    flight time of approximately one hour) and was adequately set for the approach.

215    When the crew lowered the gear the green hydraulic system malfunctioned and prevented

216    the nose gear from fully extending and locking; it could not be retracted. As a consequence, the crew

217    had to go-around and follow several procedures. With the aerodynamic drag being doubled, flight

218    endurance was halved (approximately 30 min). In their subsequent approach the crew was already in

219    a mayday situation. Upon selecting the next flap level, due to the underlying failure of the green

220    hydraulic system, the flaps or the slats (depending on the initial configuration) jammed. The high task

221    load condition for the rating experiment began at this point. As the malfunction affected the landing

222    performance of the aircraft, the crew was again requested to handle several procedures before they

223    were able to land. For further details on the technical scenario, the reader is referred to Gontar and

224    Hoermann (2014).

225    The scenario was very challenging and elicited the pilots' CRM skills on all the dimensions

226    that are usually trained and rated during recurrent training. These dimensions include

227    communication skills, leadership and teamwork, work organization as well as situation awareness

228    and decision making. Since the malfunctions that occurred were unforeseen for the pilots, they were

229    not able to prepare themselves beforehand, but had to make fast decisions, very efficient task

230    assignments, and also handle the procedures, the automation and checklists with particular

231    precision. Normally, the crews have enough time to work through their procedures step by step.

232 However, the fuel problem in our scenario forced them to work through the procedures more quickly

233 and thus communicate more concisely and more effectively. Furthermore, the success of this mission

234 was highly dependent on making the proper decisions in the right order (e.g. aborting or skipping

235 procedures or declaring an emergency due to the very low fuel level). It was expected that only

236 crews with high CRM skills would be able to complete this mission satisfactorily.

237 During the 30 flight simulator missions, we recorded audio data from the pilots and the ATC,

238 flight simulator data, as well as video data showing the two participating pilots and the cockpit

239 interior from behind. Pilot performance was rated by a flight instructor from the respective airline

240 during the missions (benchmark rating). This benchmark rating was based on the evaluation form

241 which is used in this airline (Burger, Neb, & Hoermann, 2003) and is explained below ($NTS_{item}$). We

242 found high variance within the pilots' performance ratings. The ratings included crews that were able

243 to manage the severe technical problems very quickly and very well, but also crews which were

244 unable to deal with the problems. Based on the averaged performance grade of the benchmark

245 ratings, we selected four videotapes that reflect the entire spectrum of CRM performance: *poor*,

246 *medium-low*, *medium-high*, and *outstanding.* To validate the benchmark rating, the videotape

247 selection was verified by an additional type-rating examiner.

248 **Dependent Measures**

249 The raters assessed the pilots' performance based on videotapes using three different rating

250 tools: two NTS rating tools – one on a dimension basis, one on an item basis (Burger et al., 2003) –

251 and the *Line Operations Safety Audit* (LOSA) rating tool (Klinect, Murray, Merritt, & Helmreich, 2003).

252 Examples of the content for each tool are shown in Figure 1.

253

254

255 Insert Figure 1 around here

256

257

258    **NTS rating tool on dimension basis (NTS$_{dim}$).** The raters assessed each individual pilot's

259    performance using a five-point scale ranging from *poor* (1) to *outstanding* (5), see Figure 1 at the top.

260    The following dimensions were addressed: *communication, leadership & teamwork, work*

261    *organization,* and *situation awareness & decision making* (Burger et al., 2003)*. Communication* and

262    *leadership & teamwork* are regarded as social aspects; *work organization* and *situation awareness &*

263    *decision making* are regarded as cognitive aspects (Hoermann & Neb, 2004). This rating method

264    requires the instructor to assess the pilots' performance globally across the whole videotape, but is

265    not based on single items (Brannick et al., 2002). The raters themselves have to relate specific crew

266    behaviors to the various NTS$_{dim}$ dimensions. Such rating methods require a higher degree of

267    abstraction and are expected to be more subjective and thereby less reliable than directly observable

268    behaviors (Brannick et al., 2002).

269    **NTS rating tool on item basis (NTS$_{item}$).** The raters assessed each individual pilot's

270    performance on 40 items which reflect the same four dimensions as NTS$_{dim}$, but support the rater

271    with more specific items (Burger et al., 2003), see the middle of Figure 1. The 40 items represent the

272    dimensions *communication* (10 items)*, leadership & teamwork* (15 items)*, work organization* (8

273    items)*,* and *situation awareness & decision making* (7 items), and were rated on the same five-point

274    scale, ranging from *poor* to *outstanding*, as NTS$_{dim}$. The mean value of all items of a dimension was

275    calculated to obtain a value comparable to NTS$_{dim}$ but based on items. The items of this rating tool

276    were known to the raters and are regularly used in their airline's training. It is based on the NOTECHS

277    method (Flin et al., 2003) and was adapted to the company's culture and CRM philosophy by a

278    working group comprised of subject matter experts, such as training and check pilots, aviation

279    psychologists, and human factors specialists (Burger et al., 2003). This work was influenced by the

280    results of a safety survey that the airline conducted. The purpose of this study was to analyze safety-

281    relevant events from the preceding five years. Based on this survey, Burger et al. (2003) were able to

282    identify specific factors that contributed to the events and translated them to markers. A content

283    analysis was performed to ensure that all NOTECHS markers were covered by the newly developed

284    system.

285    **Line Operations Safety Audit.** The raters assessed the pilots' performance on four

286    dimensions using 13 items that represent *planning behavioral markers* (4 items)*, execution*

287    *behavioral markers* (4 items)*, review / modify behavioral markers* (3 items)*,* and *overall behavioral*

288    *markers* (2 items). Ratings were obtained on a four-point scale from *poor* (1) to *outstanding* (4), see

289    Figure 1 at the bottom. The rating of a dimension is given by the mean of all its item values. As

290    Haeusler, Klampfer, Amacher, and Naef (2004) demonstrated, dimensions of LOSA strongly correlate

291    with dimensions of the NOTECHS system, which was the basis of the $NTS_{item}$ system used here

292    (Burger et al., 2003). In addition, LOSA incorporates aspects of technical skills as well. Since technical

293    aspects such as *automation handling* are more overt and observable, we expect higher inter-rater

294    reliability for the LOSA rating. The LOSA rating tool was sent to the raters two weeks before the

295    rating experiment, but they were not familiar with it. We used the *LOSA Descent / Approach / Land*

296    sheet (International Civil Aviation Organization, 2002, p. A-8; Klinect et al., 2003). In contrast to the

297    NTS ratings, both crew members were rated together as a team.

298    In addition to the five-point scale of both NTS tools, we derived a dichotomous pass/fail scale

299    by assigning the lower two scale points to *fail* and the upper three scale points to *pass.* Regarding the

300    LOSA rating, we assigned the lowest scale point to *fail* and the upper three scale points to *pass*; if one

301    item was rated as failed, the whole rating dimension was deemed unsatisfactory.

**Instructions and Procedure**

302

303    Two weeks before the rating experiment was conducted, all raters were informed about the

304    upcoming assessment. They were briefed about the technical details of the scenario, such as the

305    expected route, weather conditions, aircraft configuration, malfunctions, etc. Furthermore, the

306    raters received a copy of the three rating forms they would have to use (see Figure 1 with examples

307    of the content). While the $NTS_{item}$ rating tool was already known to the instructor pilots from their

308    current training practice, the $NTS_{dim}$ and the LOSA tools had not been used by the airline before.

309    Immediately before the rating experiment began, we once again explained the rating tools

310    and the whole scenario to the raters. The raters were explicitly advised to leave items blank if they

311    did not observe a corresponding behavior. In addition, the raters were instructed not to talk to each

312    other and were told that they were not allowed to page back in the rating sheets – neither during the

313    rating process itself, nor between rating the different crews. The sheets were then handed out.

314    Following these instructions, the four videotapes with a duration of *M* = 5.98 (*SD* = 1.42)

315    minutes each were presented one by one. Rewinding or repeating was not an option; however,

316    raters could take notes. The videotapes were presented in the following order of crew performance:

317    (1) *medium-high*, (2) *medium-low*, (3) *outstanding*, and (4) *poor*. In order to minimize sequence

318    effects, the medium-performing crews (1; 2) were presented first. Between the presentations, the

319    instructors rated the pilots' skills; all raters had as much time as they wanted to complete their

320    ratings. This took approximately 20 minutes after each scenario. The videotape began with a map

321    showing the current location of the aircraft, speed, heading, and the remaining fuel on board for a

322    duration of 30 seconds. The actual flight scenario was then presented, starting exactly 30 seconds

323    before the second malfunction occurred (flaps or slats jammed) and the high task load condition

324    began.

325    **Analysis**

326        For each of the dependent measures, we calculated $r_{wg}$ (cf. James, Demaree, & Wolf, 1984)

327    to assess within-group agreement for the five-point and four-point scales, and on the pass/fail level.

328    "The technique [of $r_{wg}$] was cast as a heuristic form of interrater reliability…" (James, Demaree, &

329    Wolf, 1993, p. 306) and sees total variance (in contrast to classical test theory) as being related to the

330    rater (cf. Liao, Hunt, & Chen, 2010). Values of $r_{wg}$ were calculated for each rating dimension of the

331    respective tool and for each performance level of the crew. This allowed us to identify the

332    *measurement*-related and *target*-related influences on inter-rater reliability. In addition, for $NTS_{dim}$

333    and $NTS_{item}$, $r_{wg}$ was calculated separately for the CPTs' and the FOs' performance ratings. As the crew

334    was assessed as one team by LOSA, a comparison between the crew members was not possible. The

335    threshold value for acceptable within-group agreement was set to .70 (Nunnally & Bernstein, 1994),

336    so that agreement equal to or above .70 is interpreted as *acceptable*, and values below .70 are

337    interpreted as *not acceptable* agreement; for an in-depth discussion about this commonly used

338    criterion, see Harvey and Hollander (2004).

339        Intraclass correlation coefficients *ICC(3)* for single measures were calculated using a two-way

340    mixed model for each of the dependent measures to assess inter-rater reliability (cf. Shrout & Fleiss,

341    1979). The factor *rater* was determined as fixed since the raters were preselected for the workshop.

342    In contrast to most of the studies in the medical domain, average measures of *ICC(3)* do not seem

343    appropriate here, since the reliability of one single rater, and thus the single measure, is decisive.

344    This is due to the fact that only one instructor pilot rates the crew's performance in real training and

345    examination flights. *ICC(3)* analyses were conducted at the five-point and four-point scale level only

346    (not on the derived dichotomous level). ICCs were calculated for all the different rating dimensions of

347    the two NTS measurements and the LOSA measurement.

348         Although ICCs can be calculated for true dichotomous data, a subsequently derived

349    dichotomous level from a higher level scale (as shown here) would require tetrachoric correlation

350    coefficients (Wirtz & Caspar, 2002). Based on our data set, it was not possible to calculate tetrachoric

351    correlation coefficients due to missing data and the resultant singularities. According to Cicchetti

352    (1994) who subdivided the recommendation of Landis and Koch (1977), the *ICC(3)* values are

353    interpreted as follows: Values below .40 represent *poor* clinical significance, values between .40 and

354    .59 represent *fair*, and values between .60 and .75 *good* clinical significance. Values greater than .75

355    are considered *excellent* clinical significance (Cicchetti, 1994; see also Fleiss, Levin, and Paik (2003).

356    When calculating mean values of ICCs, *Fisher z' transformation* (Fisher, 1925) was used. With respect

357    to the ICCs, the Spearman-Brown prophecy formula (cf. Lienert & Raatz, 1998) was used to calculate

358    the minimum number of raters that is required to achieve a specific level of reliability (e.g., .60 for

359    good clinical significance according to Cicchetti, 1994).

360         According to the model assumptions for *ICC(3)* made by Shrout and Fleiss (1979), we used

361    Shapiro-Wilk tests to analyze for normal distribution as suggested by Thode (2002) and Razali and

362    Wah (2011). Analyses showed that the sample data were non-normally distributed ($p < .05$), except

363    for the LOSA dimensions *execution behavioral markers, W*(143) = .99, *p* = .39, and *overall behavioral*

364    *markers, W*(143) = .98, *p* = .06. Based on the visual inspection of the plots, we concluded that the

365    significant results in the tests were rather due to the large sample size than due to meaningful

366    deviations from the normal distribution (Field, 2009). The residuals showed non-normal distributions

367    ($p < .05$) as well, except for the $NTS_{item}$ dimensions *communication, W*(288) = .99, *p* = .45, *leadership*

368    *& teamwork D*(288) = .99, *p* = .20*, situation awareness & decision making W*(288) = .99, *p* = .72, and

369    for the LOSA dimension *review / modify, W*(136) = .99, *p* = .24. Since the *analysis of variance,* which

370    corresponds to the *ICC(3)* model, is robust against violations of normal distribution (Schmider,

371    Ziegler, Danay, Beyer, & Buehner, 2010), we did not anticipate problems in using ICCs for these data.

372                                        **Results**

373    **Results Regarding the $NTS_{dim}$ Tool**

374    Table 1 illustrates the results with respect to the analysis of NTS$_{dim}$. Regarding the within-
375    group agreement $r_{wg}$ on the five-point scale, the results showed acceptable agreement for the CPT
376    (.79 on average) and for the FO (.74 on average) in the videotape that showed crew members with
377    *outstanding* performance. The performance of all other pilots was rated with an agreement lower
378    than .70 and was therefore not acceptable. When looking at the average of the different rating
379    dimensions, it can be seen that the agreement across all videotapes was below the .70 threshold for
380    every dimension. In three out of four videotapes, the FOs' performance was rated in higher
381    agreement than the CPTs'; the average agreement across all videotapes (.57) was not acceptable. It
382    can be concluded that the agreement of raters depended on the level of performance that was
383    exhibited by the pilots.

384    *ICC(3)* for inter-rater reliability was found to be poor for the dimensions *communication*
385    (.12), *leadership & teamwork* (.28), and *work organization* (.34) of NTS$_{dim}$. Only ratings for *situation
386    awareness & decision making* (.45) represented fair reliability. Inter-rater reliability of social aspects
387    (*communication* and *leadership & teamwork*) was lower than for cognitive aspects (*work
388    organization* and *situation awareness & decision making*). The average inter-rater reliability was poor
389    (.30). In order to reach a good level of reliability with respect to the ICCs (.60), the Spearman-Brown
390    prophecy formula revealed that on average, four raters would be required to assess pilots' non-
391    technical skills on the dimensional level for such scenarios.

392

393 ─────────────────────────────────────────────

394                              Insert Table 1 around here

395

396 ─────────────────────────────────────────────

397    Looking at the within-group agreement results from the derived pass/fail scale (see Table 1,
398    bottom), once again the *outstanding* performing crew was the only one which was represented by
399    acceptable ratings for the CPT (.97) and for the FO (.92). While the ratings for this crew include six
400    ratings that were in perfect agreement (1.0), the *medium-high* and *poor* performing crews included
401    ratings that showed no agreement (0.0). The mean $r_{wg}$ of the rating dimensions were all below the

402  required minimum of .70. Comparing the pass/fail scale with the five-point scale, the *outstanding*

403  performing crew was rated in higher agreement on the pass/fail scale (.97 / .92 vs. .79 / .74). The

404  opposite was true for the *medium-high* and *poor* performing crews. The average agreement was

405  lower on the dichotomous pass/fail scale (.46) than on the five-point scale (.57).

406  **Results Regarding the NTS$_{item}$ Tool**

407  Within-group agreement $r_{wg}$ showed acceptable ratings for the *poor* (.71 / .71) and

408  *outstanding* (.77 / .80) performing crews as well as for the FO of the *medium-low* (.78) performing

409  crew (see Table 2, top). All dimensions of NTS$_{item}$ showed acceptable agreement, although they were

410  close to or at the threshold of .70. The trend indicated in the results of NTS$_{dim}$, i.e. that the FOs'

411  performance was rated as slightly more in agreement than the CPTs' was seen here as well. A

412  comparison of the mean of the dimensions shows that NTS$_{item}$ ratings were in higher agreement than

413  NTS$_{dim}$ ratings on the five-point scale.

414  As already measured for NTS$_{dim}$, *ICC(3)* reliability was fair for *situation awareness & decision*

415  *making* (.48), but poor for all the other dimensions; the trend that social aspects are rated less

416  reliably than cognitive aspects was reflected in these results as well. Comparing the NTS$_{item}$ and the

417  NTS$_{dim}$, it can be seen that the inter-rater reliability for *communication* was higher for NTS$_{item}$ than for

418  NTS$_{dim}$. With respect to the ICCs, the Spearman-Brown prophecy formula found that on average,

419  three raters would be sufficient to assess pilots' non-technical skills on this five-point scale with good

420  (.60) reliability.

421

422

423  Insert Table 2 around here

424

425

426  The pass/fail scale showed unacceptable agreement and thus lower agreement than the five-

427  point scale in every single value (see Table 2, bottom). In contrast to the five-point scale, the average

428  agreement for the FOs' performance ratings was lower than for the CPTs' performance ratings. The

429  average agreement using the NTS$_{dim}$ tool (.46) was higher than the average NTS$_{item}$ agreement (.19)

430  on the pass/fail scale.

431  **Results Regarding the LOSA Tool**

432      Results regarding the LOSA rating showed acceptable within-group agreement for the

433  *planning* (.74) and *execution* (.76) dimensions, but agreement below the defined .70 threshold for

434  the *review / modify* (.63) and *overall* (.61) dimensions (see Table 3, top). Rating for the *outstanding*

435  (.74) and *medium-low* (.71) performing crews showed acceptable agreement on average. Agreement

436  for *poor* (.68) and *medium-high* (.61) performing crews was slightly below the threshold. Inter-rater

437  reliability was fair for the *planning* (.47) and *execution* (.43) dimension, but was poor for the *review /*

438  *modify* (.25), and *overall* (.30) dimensions. Although the raters had not used this rating tool in their

439  training before, the average reliability of LOSA (.37) was slightly higher than for NTS$_{dim}$ (.30) and

440  NTS$_{item}$ (.35). With respect to the ICCs, the Spearman-Brown prophecy formula found that on

441  average, three raters would be required to assess pilots' skills on the LOSA scale with good (.60)

442  reliability.

443

444  _____

445                        Insert Table 3 around here

446  _____

447

448      Agreement on the pass/fail scale (see Table 3, bottom) was high (.92) for the *outstanding*

449  performing crew. Agreement for the lower performing crews was below the acceptable threshold.

450  On average, the rating dimensions did not exceed the acceptable threshold. As for NTS$_{dim}$ and NTS$_{item}$,

451  the average agreement for *medium-high*, *medium-low*, and *poor* performing crews was lower for the

452  pass/fail scale than for the four-point scale. Both rating tools, which had not been used by the raters

453  before (NTS$_{dim}$ and LOSA), showed higher agreement for the pass/fail scale than for the five-

454  point/four-point scale when rating the *outstanding* performing crew.

455                                **Discussion**

456     The discussion is divided into several parts that correspond to the *measurement* and *target*

457     themes introduced in this paper. It concludes by pointing out some limitations of the study.


***Measurement*: Different Rating Tools (NTS$_{dim}$, NTS$_{item}$, LOSA) and Familiarity**

459     The results showed that neither of the rating tools used here achieved the necessary

460     standard for sufficient inter-rater reliability on average across all dimensions and performance levels.

461     Although the raters were not trained with the LOSA sheet, the reliability was roughly the same for

462     this rating tool at the four-point scale as compared to the other rating tools. The average agreement

463     with LOSA on the pass/fail scale was even better than agreement with the rating tool known from

464     training (NTS$_{item}$). This means on one hand that familiarity with the rating tool alone does not

465     necessarily lead to higher inter-rater reliability. On the other hand these results indicate that more

466     precisely formulated items (LOSA) can outweigh the potential familiarity advantages (NTS$_{item}$). In this

467     context it has to be mentioned that NOTECHS (which was the basis of NTS$_{item}$) was not intended to be

468     used on a pass/fail level unless a rating could be linked to technical consequences (Flin et al., 2003).


469     Another reason for the similar inter-rater reliability of the LOSA tool could be that the crew is

470     rated as one team (LOSA) and not as two single pilots (as for NTS$_{dim}$ and NTS$_{item}$). Perhaps it is more

471     difficult for raters to assign separate performance contributions to the two crew members, since

472     interaction, which is the basis for most of the NTS dimensions introduced here, is the result of a

473     collaboration of at least two persons. O'Connor et al. (2002) came to a similar conclusion when

474     addressing inter-rater reliability differences between scenarios. Another aspect could be that LOSA

475     also incorporates technical performance aspects such as *automation handling* that are easier to

476     observe. Finally, LOSA only uses 13 items in contrast to NTS$_{item}$, which uses 40 items; raters could

477     have lost interest in thoroughly considering the item definitions before assigning a score. In our

478     study, they had to go through all ratings eight times, which could have led to a *checking-boxes*

479     response style. Although the agreement for NTS$_{item}$ on the five-point scale is acceptable on average,

480     the assessment on the pass/fail scale is what an examination ultimately depends on.


***Measurement*: Scale Levels and Degree of Differentiation**

482     When comparing the two scales, which represent different levels of differentiation, all the

483     rating tools showed lower agreement on the derived pass/fail scale than on the five-point/four-point

484     scale (NTS$_{dim}$: .46 vs. .57; NTS$_{item}$: .19 vs. .72; LOSA: .43 vs. .68) on average. It seems that examiners

17

485 can give reliable feedback in general (e.g. pilot A was better than pilot B), but are in less agreement

486 with respect to the level of pass and fail (e.g. pilot A passed, but pilot B failed). O'Connor et al. (2002)

487 found similar results for two of their eight scenarios. This issue is relevant in particular because this

488 threshold between pass and fail is what counts most for the individual pilot. This finding confirms

489 earlier concerns that NTS ratings alone should not be used to pass or fail a crew member unless

490 safety consequences are directly involved.

491 *Measurement*: **Differences Between the Rating Dimensions**

492 The results showed that inter-rater reliability was dependent on the dimension being rated.

493 Comparing the different dimensions, both $NTS_{dim}$ and $NTS_{item}$ showed lower inter-rater reliability for

494 the social aspects than for the cognitive aspects. When we assessed the reliability of the pilots' self,

495 peer, and supervisor ratings, we also found less agreement for the social aspects than for the

496 cognitive aspects using the $NTS_{item}$ tool (cf. Gontar & Hoermann, 2014). In this earlier study, the

497 entire mission was rated by all 60 pilots who took part in the simulator experiment. When the whole

498 mission was rated, inter-rater reliability for the social dimensions was even lower than in the present

499 study. The opposite was true for the cognitive aspects, which led to higher inter-rater reliabilities

500 when rating the entire mission and when the rater was directly involved and operating the simulator.

501 In contrast to the findings presented here, Yule et al. (2008) found that for surgeons, aspects of

502 communication, leadership and teamwork (which correspond to our social skills) were rated more

503 reliably than aspects of task management, decision making and situation awareness (which

504 correspond to our cognitive skills). One reason for this different finding could be that the videotapes

505 we selected for this rater study featured a strong emphasis on aspects of decision making and

506 situation awareness. This is because the success in this scenario mainly depended on the appropriate

507 decision making by the crew. This aspect might also be the reason why the LOSA dimension of

508 *Planning* is rated with slightly higher inter-rater reliability than the other LOSA dimensions.

509 *Target*: **Crews Representing Different Levels of Performance**

510 Based on the work from Yule et al. (2009), we expected that the most extreme performance

511 characteristics, such as the *outstanding* and *poor* performing crews, would be rated with higher

512 agreement than the *medium* performing crews. What we found was that only the *outstanding*

513 performing crew was rated with acceptable agreement on average. That the *poor* performing crew

514    was rated with lower agreement on the pass/fail scale than the *outstanding* performance is even

515    more surprising, because it was rated directly after the latter. This may indicate that the raters were

516    not subject to sequence effects; had that been the case, they would have consistently rated the *poor*

517    performing crew as very poor. An explanation could be that the raters are seldom faced with *poor*

518    performing crews so that they struggle to assess them in high agreement.

*Target***: Differences Between the Ratings for Captains and First Officers**

520    The agreement using the two NTS rating tools showed that the FOs are rated slightly more in

521    agreement on the five-point scale than the CPTs, except using the NTS$_{dim}$ tool for the *outstanding*

522    performing crew. O'Connor et al. (2002) compared the deviations of ratings for CPTs and FOs to a

523    reference rating. They found similar results with the ratings for the FOs' performance showing less

524    deviation from the reference ratings than the ratings for the CPTs' performance. It seems that the

525    raters can assess FOs' performance more accurately than the CPTs'. As all the raters were CPTs

526    themselves, they were used to flying together with FOs more often than with CPTs. This daily

527    experience could have provided them with better framing conditions to compare the behaviors of

528    the FOs. Another aspect is that in examination flights for example, a good pilot can

529    disproportionately influence the team's performance and thus compensate for the effects of a more

530    poorly performing crew member. If the good pilot continuously supports his crew member, the

531    performance by the poorly performing pilot might be hidden. So even when both crew members are

532    rated independently (NTS$_{dim}$ and NTS$_{item}$), it might be hard to identify the differences between two

533    pilots' non-technical skills.

**Strengths and Weaknesses of this Study**

535    One weakness of this study might be that the raters did not have enough time to actually

536    observe every behavior, because they only observed a 6-minute sequence of the whole mission.

537    Moreover, they may have assessed aspects which they did not observe instead of leaving the

538    respective items blank or marking it as *not observable* (as they were instructed to do). Such behavior

539    has already been reported by O'Connor et al. (2002). Another effect of the rather short duration and

540    the content of the videotape might have been that the raters did not have the opportunity to

541    observe the pilots' behavior during normal operation flight phases. This would have been the case in

542    examination flights, which usually start with normal operations before the crews are exposed to

543     critical situations. Nevertheless, it is especially the performance in unforeseen and abnormal events

544     that determines the success of a mission.

545     In contrast, a strength of this study is the large sample of non-volunteer raters which

546     represents the most experienced instructors in the entire airline fleet. We used a clean environment,

547     meaning the raters did not have any parallel tasks such as operating the simulator or acting as air

548     traffic controllers, which could have led to high rater workload during the assessment (Deaton et al.,

549     2007; Seamster, Hamman, & Edens, 1995). Furthermore, all videotapes contained the same task with

550     different performance levels, as would be expected from daily practice.

551     **Conclusion**

552     The results of this inter-rater reliability study show that the *measurement* as well as the

553     *target* influence inter-rater reliability. We were able to show these effects while keeping other

554     influences by the *rater* and the *scenario & task* constant. On the other hand, we demonstrated that

555     inter-rater reliability is still an unsolved issue even within a group of highly experienced instructors

556     when assessing the non-technical skills of pilots. In Europe, current regulatory material by the

557     European Aviation Safety Agency (2011) states that the practical training of instructors should

558     "include the development of specific instructor skills, particularly in the area of teaching and

559     assessing threat and error management and CRM" (European Aviation Safety Agency, 2011, FCL.920,

560     p. 282). In particular, instructors are required to observe and assess CRM behaviors in order to

561     provide constructive feedback to both the pilots and to the training department (European Aviation

562     Safety Agency, 2011). All these requirements assume that such ratings are based on reliable

563     observations. According to our findings, we strongly recommend incorporating specific inter-rater

564     reliability exercises into trainer standardization and assessment trainings. Therefore, it would be

565     beneficial to describe more precise anchors for desired and undesired behaviors on all observed CRM

566     dimensions. Based on our findings we recommend caution when using NTS-ratings on a pass/fail

567     level. In line with the second NOTECHS principle it should be emphasized that in order to fail a pilot

568     in an examination flight due to non-technical skills, "flight safety must be actually (or potentially)

569     comprised [, which] requires a related objective technical consequence" (Flin et al., 2003, p. 109).

570     We agree with the opinion that mission-specific CRM evaluation tools and other objectifying

571     resources, such as shown by Brannick et al. (2002), would lead to higher inter-rater reliability in

572     training. Deaton et al. (2007) for example developed a tool which supports the instructors when

573  rating pilots' performance in specific scenarios. This tool will alert the instructor when it detects

574  events in training scenarios that are important for the rating. These authors could show that such a

575  supporting tool leads to more differentiated and more accurate ratings (Deaton et al., 2007). Such an

576  approach seems to be very promising. A goal of future research should be to further elaborate and

577  extend the usage of such techniques with the goal of providing the instructors with more reliable

578  information for their assessments. In parallel, regulatory authorities should explicitly advise airlines

579  and flight training organizations to address and demonstrate sufficient inter-rater reliability among

580  their instructor pilots when utilizing their performance evaluation tools.

581

582  Acknowledgements

Gontar, P. & Hoermann, H.-J. (2015). Inter-rater reliability at the top end: Measures of pilots' non-technical performance. The International Journal of Aviation Psychology, 25(3/4), 171-190.

588                                        References

589   Arora, S., Miskovic, D., Hull, L., Moorthy, K., Aggarwal, R., Johannsson, H., . . . Sevdalis, N. (2011). Self

590          vs. expert assessment of technical and non-technical skills in high fidelity simulation. *The*

591          *American Journal of Surgery, 202*(4), 500–506. doi:10.1016/j.amjsurg.2011.01.024

592   Beard, J. D., Marriott, J., Purdie, H., & Crossley, J. (2011). Assessing the surgical skills of trainees in the

593          operating theatre: A prospective observational study of the methodology. *Health Technology*

594          *Assessment, 15*(1), 1–162. doi:10.3310/hta15010

595   Brannick, M. T., & Prince, C. (1991). *Assessment of aircrew rating from within and between scenarios*

596          *(Tech. Rep. No. DAAL0–3–86–D–001).* Orlando, FL.

597   Brannick, M. T., Prince, C., & Salas, E. (2002). The reliability of instructor evaluations of crew

598          performance: Good news and not so good news. *The International Journal of Aviation*

599          *Psychology, 12*(3), 241–261. doi:10.1207/S15327108IJAP1203_4

600   Burger, K.-H. (Ed.). (1999). *Basic competence for optimum performance*. Frankfurt.

601   Burger, K.-H., Neb, H., & Hoermann, H.-J. (2003). Lufthansa's new basic performance of flight crew

602          concept - A competence based marker system for defining pilots performance profile. In R. S.

603          Jensen (Ed.), *Proceedings of the 12th International Symposium on Aviation Psychology*

604          (pp. 172–175). Dayton, OH: Wright State University Press.

605   Cicchetti, D. V. (1994). Guidelines, criteria, and rules of thumb for evaluating normed and

606          standardized assessment instruments in psychology. *Psychological Assessment, 6*(4), 284–

607          290. doi:10.1037/1040-3590.6.4.284

608   Cooper, S. J., Endacott, R., & Cant, R. P. (2010). Measuring non-technical skills in medical emergency

609          care: A review of assessment measures. *Open Access Emergency Medicine, 2,* 7–16.

610          doi:10.2147/OAEM.S6693

611   Cooper, G. E., White, M. D., & Lauber, J. K. (Eds.) (1980). *Resource management on the*

612          *flightdeck: Proceedings of a NASA/Industry Workshop.* NASA Conference Publication 2120,

613          Moffett Field, CA.

614   Deaton, J. E., Bell, B., Fowlkes, J., Bowers, C., Jentsch, F., & Bell, M. A. (2007). Enhancing team

615          training and performance with automated performance assessment tools. *The International*

616          *Journal of Aviation Psychology, 17*(4), 317–331. doi:10.1080/10508410701527662

617   Dedy, N. J., Szasz, P., Louridas, M., Bonrath, E. M., Husslein, H., & Grantcharov, T. P. (2015). Objective

618          structured assessment of nontechnical skills: Reliability of a global rating scale for the in-

619      training assessment in the operating room. *Surgery, 157*(6), 1002–1013*.*

620      doi:10.1016/j.surg.2014.12.023

621  European Aviation Safety Agency. (2011). *Annex to ED decision 2011/016/R: Acceptable means of*

622      *compliance and guidance material to part-FCL.*

623  European Aviation Safety Agency. (2014). *Annex to ED Decision 2014/022/R*.

624  European Commission. (2011). *Commission regulation (EU) No 1178/2011 of 3 November 2011:*

625      *Laying down technical requirements and administrative procedures related to civil aviation*

626      *aircrew pursuant to regulation (EC) No 216/2008 of the European parliament and of the*

627      *council.*

628  Federal Aviation Administration (2004). *Advisory Circular No 120-51E: Crew resource management*

629      *training*. Washington, DC: U.S. Department of Transportation.

630  Field, A. P. (2009). *Discovering statistics using SPSS*. Los Angeles, CA: SAGE Publications.

631  Fisher, R. A. (1925). *Statistical methods for research workers*. Edinburgh: Oliver and Boyd.

632  Fleiss, J. L., Levin, B., & Paik, M. C. (2003). *Statistical methods for rates and proportions*. *Wiley Series*

633      *in probability and statistics*. Hoboken, NJ: J. Wiley.

634  Fletcher, G., Flin, R., McGeorge, P., Glavin, R., Maran, N., & Patey, R. (2003). Anaesthetists' Non-

635      Technical Skills (ANTS): Evaluation of a behavioural marker system. *British Journal of*

636      *Anaesthesia, 90*(5), 580–588*.* doi:10.1093/bja/aeg112

637  Flin, R., & Martin, L. (2001). Behavioral markers for crew resource management: A review of current

638      practice. *The International Journal of Aviation Psychology, 11*(1), 95–118.

639      doi:10.1207/S15327108IJAP1101_6

640  Flin, R., Martin, L., Goeters, K.-M., Hoermann, H.-J., Amalberti, R., Valot, C., & Nijhuis, H. (2003).

641      Development of the NOTECHS (non-technical skills) system for assessing pilots' CRM skills.

642      *Human Factors and Aerospace Safety, 3*(2), 95–117.

643  Gale, T. C. E., Roberts, M. J., Sice, P. J., Langton, J. A., Patterson, F. C., Carr, A. S., . . . Davies, P. R. F.

644      (2010). Predictive validity of a selection centre testing non-technical skills for recruitment to

645      training in anaesthesia. *British Journal of Anaesthesia, 105*(5), 603–609.

646      doi:10.1093/bja/aeq228

647  Goldsmith, T. E., & Johnson, P. J. (2002). Assessing and improving evaluation of aircrew performance.

648      *The International Journal of Aviation Psychology, 12*(3), 223–240.

649      doi:10.1207/S15327108IJAP1203_3

Gontar, P. & Hoermann, H.-J. (2015). Inter-rater reliability at the top end: Measures of pilots' non-technical performance. The International Journal of Aviation Psychology, 25(3/4), 171-190.

650  Gontar, P., & Hoermann, H.-J. (2014). Flight crew performance and CRM ratings based on three
651      different perceptions. In A. Droog (Ed.), *Aviation Psychology: Facilitating change(s):*
652      *Proceedings of the 31st EAAP Conference* (pp. 310–316). Malta.
653  Gontar, P., Hoermann, H.-J., Deischl, J., & Haslbeck, A. (2014). How pilots assess their non-technical
654      performance - A flight simulator study. In N. A. Stanton, S. J. Landry, G. Di Bucchianico, & A.
655      Vallicelli (Eds.), *Advances in Human Aspects of Transportation. Part I* (pp. 119–128). Krakow.
656  Hamman, W., & Holt, R. (1997). Line operational evaluation (LOE): Air carrier scenario based
657      evaluation. In E. Smith (Ed.), *Proceedings of the Human Factors and Ergonomics Society 41st*
658      *annual Meeting* (pp. 907–911). Albuquerque, NM: Human Factors and Ergonomics Society.
659  Harvey, R. J., & Hollander, E. (2004). Benchmarking $r_{wg}$ interrater agreement indices: Let's drop the
660      .70 rule-of-thumb. In *Annual Conference of the Society for Industrial and Organizational*
661      *Psychology*. Chicago, IL.
662  Haeusler, R., Klampfer, B., Amacher, A., & Naef, W. (2004). Behavioral markers in analyzing team
663      performance of cockpit crews. In R. Dietrich & T. M. Childress (Eds.), *Group interaction in*
664      *high risk environments*. Aldershot: Ashgate.
665  Helmreich, R. L., Merritt, A. C., & Wilhelm, J. A. (1999). The evolution of crew resource management
666      training in commercial aviation. *The International Journal of Aviation Psychology, 9*(1), 19–32.
667      doi:10.1207/s15327108ijap0901_2
668  Helmreich, R. L., Klinect, J.R., Wilhelm, J. A., & Jones, S.G. (1999). *The Line/LOS Error Checklist,*
669      *Version 6.0: A checklist for human factors skills assessment, a log for off-normal events, and a*
670      *worksheet for cockpit crew error management*. Austin, TX.
671  Hoermann, H.-J. & Neb, H. (2004). From NOTECHS to LH behavior markers: An implementation case
672      study. Paper presented to the Royal Aeronautical Society - Human Factors Group Seminar on
673      Assessment & Accreditation, London April 30, 2004.
674  Holt, R. W., Hansberger, J. T., & Boehm-Davis, D. A. (2002). Improving rater calibration in aviation: A
675      case study. *The International Journal of Aviation Psychology, 12*(3), 305–330.
676      doi:10.1207/S15327108IJAP1203_7
677  International Civil Aviation Organization. (2002). Line Operations Safety Audit (LOSA): DOC 9803.
678      AN/761. Montreal.
679  James, L. R., Demaree, R. G., & Wolf, G. (1984). Estimating within-group interrater reliability with and
680      without response bias. *Journal of Applied Psychology, 69*(1), 85–98. doi:10.1037/0021-
681      9010.69.1.85

Gontar, P. & Hoermann, H.-J. (2015). Inter-rater reliability at the top end: Measures of pilots' non-technical performance. The International Journal of Aviation Psychology, 25(3/4), 171-190.

682  James, L. R., Demaree, R. G., & Wolf, G. (1993). $r_{wg}$: An assessment of within-group interrater
683       agreement. *Journal of Applied Psychology, 78*(2), 306–309. doi:10.1037/0021-9010.78.2.306
684  Joint Aviation Authorities. (2001*). Joint aviation regulations: JAR OPS 1.940, 1.945, 1.955, and 1.965.*
685       Hoofddorp, Netherlands.
686  Klinect, J. R., Murray, P., Merritt, A. C., & Helmreich, R. L. (2003). Line Operations Safety Audit (LOSA)
687       - Definition and operating characteristics. In *Proceedings of the 12th International*
688       *Symposium on Aviation Psychology* (pp. 663–668). Dayton, OH.
689  Kontogiannis, T., & Malakis, S. (2013). Strategies in coping with complexity: Development of a
690       behavioural marker system for air traffic controllers. *Safety Science, 57,* 27–34.
691       doi:10.1016/j.ssci.2013.01.014
692  Landis, R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data.
693       *Biometrics, 33*(1), 159–174.
694  Liao, S. C., Hunt, E. A., & Chen, W. (2010). Comparison between inter-rater reliability and inter-rater
695       agreement in performance assessment. *Annals Academy of Medicine Singapore, 39*(8), 613–
696       618.
697  Lienert, G. A., & Raatz, U. (1998). *Testaufbau und Testanalyse* [Test construction and test analysis].
698       Weinheim: Beltz, Psychologie Verl.-Union.
699  Mishra, A., Catchpole, K., & McCulloch, P. (2009). The Oxford NOTECHS System: Reliability and
700       validity of a tool for measuring teamwork behaviour in the operating theatre. *Quality and*
701       *Safety in Health Care, 18*(2), 104–108. doi:10.1136/qshc.2007.024760
702  Mitchell, L., Flin, R., Yule, S., Mitchell, J., Coutts, K., & Youngson, G. (2012). Evaluation of the scrub
703       practitioners' list of intraoperative non-technical skills (SPLINTS) system. *International*
704       *Journal of Nursing Studies, 49*(2), 201–211. doi:10.1016/j.ijnurstu.2011.08.012
705  Nunnally, J. C., & Bernstein, I. H. (1994). *Psychometric theory*. *McGraw-Hill series in psychology*. New
706       York: McGraw-Hill.
707  O'Connor, P., Hoermann, H.-J., Flin, R., Lodge, M., & Goeters, K.-M. (2002). Developing a method for
708       evaluating crew resource management skills: A European perspective. *The International*
709       *Journal of Aviation Psychology, 12*(3), 263–285. doi:10.1207/S15327108IJAP1203_5
710  Razali, N. M., & Wah, Y. B. (2011). Power comparisons of Shapiro-Wilk, Kolmogorov-Smirnov,
711       Lilliefors and Anderson-Darling tests. *Journal of Statistical Modeling and Analytics, 2*(1), 21–
712       33.

713  Schmider, E., Ziegler, M., Danay, E., Beyer, L., & Buehner, M. (2010). Is it really robust?

714        Reinvestigating the robustness of ANOVA against violations of the normal distribution

715        assumption. *Methodology, 6*(4), 147–151. doi:10.1027/1614-2241/a000016

716  Seamster, T., Hamman, W., & Edens, E. (1995). Specification of observable behaviors within

717        LOE/LOFT event sets. In R. S. Jensen (Ed.), *Proceedings of the 8th Symposium of Aviation*

718        *Psychology* (pp. 663–668). Columbus, OH: Ohio State University.

719  Sevdalis, N., Davis, R., Koutantji, M., Undre, S., Darzi, A., & Vincent, C. A. (2008). Reliability of a

720        revised NOTECHS scale for use in surgical teams. *The American Journal of Surgery, 196*(2),

721        184–190. doi:10.1016/j.amjsurg.2007.08.070

722  Shrout, P. E. Fleiss, J. L. (1979). Intraclass Correlations: Uses in assessing rater reliability.

723        *Psychological Bulletin, 86*(2), 420–428.

724  Thode, H. C. (2002). *Testing for normality*. New York: CRC Press.

725  Weber, D. E., Roth, W.-M., Mavin, T. J., & Dekker, S. W. (2013). Should we pursue inter-rater

726        reliability or diversity? An empirical study of pilot performance assessment. *Aviation in Focus*

727        *– Journal of Aeronautical Sciences*, *4*(2), 34–58.

728  Weber, D. E., Mavin, T. J., Roth, W. M., Henriqson, E., & Dekker, S. W. A. (2014). Exploring the use of

729        categories in the assessment of airline pilots' performance as a potential source of

730        examiners' disagreement. *Journal of Cognitive Engineering and Decision Making, 8*(3), 248–

731        264. doi:10.1177/1555343414532813

732  Williams, D., Holt, R., & Boehm-Davis, D. (1997). Training for inter-rater reliability: Baselines and

733        benchmarks. In R. S. Jensen (Ed.), *Proceedings of the 9th Symposium on Aviation Psychology*

734        (pp. 514–519). Columbus, OH.

735  Wirtz, M., & Caspar, F. (2002). *Beurteilerübereinstimmung und Beurteilerreliabilität: Methoden zur*

736        *Bestimmung und Verbesserung der Zuverlässigkeit von Einschätzungen mittels*

737        *Kategoriensystemen und Ratingskalen* [Rater agreement and rater reliability: Methods to

738        measure and to improve the reliability of assessments using categorical data and rating

739        scales]. Goettingen: Hogrefe.

740  Yule, S., Flin, R., Maran, N., Rowley, D., Youngson, G., & Paterson-Brown, S. (2008). Surgeons' non-

741        technical skills in the operating room: Reliability testing of the NOTSS behavior rating system.

742        *World Journal of Surgery, 32*(4), 548–556. doi:10.1007/s00268-007-9320-z

743  Yule, S., Rowley, D., Flin, R., Maran, N., Youngson, G., Duncan, J., & Paterson-Brown, S. (2009).

744        Experience matters: Comparing novice and expert ratings of non-technical skills using the

Gontar, P. & Hoermann, H.-J. (2015). Inter-rater reliability at the top end: Measures of pilots' non-technical performance. The International Journal of Aviation Psychology, 25(3/4), 171-190.

745        NOTSS system. *ANZ Journal of Surgery, 79*(3), 154–160. doi:10.1111/j.1445-

746        2197.2008.04833.x

747

$NTS_{dim}$

| CPT | | | | | CRM-rating dimension | FO | | | | |
|-----|---|---|---|----|----------------------|-----|---|---|---|----|
| ++ | + | o | - | -- | Communication | ++ | + | o | - | - - |

*"Communication includes information transfer and social aspects. The crew members share their information, and assure reception and understanding. Suggestions of other crew members are considered, even if one does not agree. Ambiguities and uncertainties are announced."* (Burger, 1999, p. 14)

$NTS_{item}$

| CPT | | | | | Communication | FO | | | | |
|-----|---|---|---|----|----------------------|-----|---|---|---|----|
| ++ | + | o | - | -- | announce ambiguities | ++ | + | o | - | - - |

LOSA (from International Civil Aviation Organization, 2002, p. A-8; Klinect et al., 2003)

| 1 = Poor | 2 = Marginal | 3 = Good | 4 = Outstanding |
|----------|--------------|----------|-----------------|
| Observed performance had safety implications | Observed performance was barely adequate | Observed performance was effective | Observed performance was truly noteworthy |
| Overall behavioral markers | | | Crew Rating |
| Communication environment | Environment for open communication was established and maintained | Good cross talk – flow of information was fluid, clear, and direct | |

748

749 Figure 1. Examples of content for the different rating tools that were used for the inter-rater

750 reliability study. From top to bottom: $NTS_{dim}$ (see Burger et al., 1999), $NTS_{item}$ (see Burger et al.,

751 2003), LOSA (see International Civil Aviation Organization, 2002, p. A-8; Klinect et al., 2003). For

752 $NTS_{dim}$, the definitions of the items were not specified on the rating tool, but were included in the

753 airline's training material the raters had.

754

Gontar, P. & Hoermann, H.-J. (2015). Inter-rater reliability at the top end: Measures of pilots' non-technical performance. The International Journal of Aviation Psychology, 25(3/4), 171-190.

755 Table 1. $r_{wg}$ and *ICC(3)* of the NTS$_{dim}$ ratings as a function of the performance level shown in the

756 scenario and crew position.

757

| NTS$_{dim}$ (CPT/FO) | Performance Level | | | | $r_{wg}$ mean | *ICC(3)* |
|---|---|---|---|---|---|---|
| | Out-standing | Medium-high | Medium-low | Poor | | |
| **Five-point scale** | | | | | | |
| Communication | .77 / .62 | .33 / .62 | .53 / .67 | .42 / .53 | .56 | .12 |
| Leadership & Teamwork | .84 / .83 | .42 / .62 | .25 / .66 | .62 / .63 | .61 | .28 |
| Work Organization | .75 / .77 | .37 / .57 | .56 / .66 | .36 / .48 | .57 | .34 |
| Situation Awareness & Decision Making | .81 / .73 | .47 / .29 | .50 / .53 | .44 / .55 | .54 | .45 |
| Mean | .79 / .74 | .40 / .53 | .46 / .63 | .46 / .55 | .57 | .30 |
| **Pass/fail scale** | | | | | | |
| Communication | .89 / .69 | .10 / .60 | .51 / .51 | .29 / .51 | .51 | - |
| Leadership & Teamwork | 1 / 1 | .00 / .59 | .50 / .59 | .29 / .58 | .57 | - |
| Work Organization | 1 / 1 | .00 / .14 | .59 / .43 | .00 / .13 | .41 | - |
| Situation Awareness & Decision Making | 1 / 1 | .10 / .00 | .59 / .17 | .00 / .00 | .36 | - |
| Mean | .97 / .92 | .05 / .33 | .55 / .43 | .14 / .30 | .46 | - |

758

759 *Note*. In addition to the five-point scale (top), we derived a dichotomous pass/fail scale by assigning

760 the lower two scale-points to *fail* and the upper three scale-points to *pass* (bottom). The mean ICC

761 value was calculated using Fisher z' transformation.

762  Table 2. $r_{wg}$ and *ICC(3)* of the NTS$_{item}$ ratings as a function of the performance level shown in the

763  scenario and crew position.

764

| NTS$_{item}$ (CPT/FO) | Performance Level | | | | $r_{wg}$ mean | *ICC(3)* |
|---|---|---|---|---|---|---|
| | Out-standing | Medium-high | Medium-low | Poor | | |
| **Five-point scale** | | | | | | |
| Communication | .76 / .77 | .54 / .68 | .69 / .76 | .74 / .75 | .71 | .22 |
| Leadership & Teamwork | .79 / .84 | .73 / .71 | .75 / .79 | .76 / .71 | .76 | .32 |
| Work Organization | .76 / .79 | .52 / .64 | .69 / .79 | .64 / .75 | .70 | .37 |
| Situation Awareness & Decision Making | .78 / .81 | .60 / .63 | .64 / .79 | .69 / .62 | .70 | .48 |
| Mean | .77 / .80 | .60 / .67 | .69 / .78 | .71 / .71 | .72 | .35 |
| **Pass/fail scale** | | | | | | |
| Communication | .43 / .13 | .24 / .10 | .00 / .00 | .13 / .05 | .13 | - |
| Leadership & Teamwork | .17 / .43 | .30 / .10 | .00 / .00 | .36 / .00 | .17 | - |
| Work Organization | .59 / .59 | .30 / .09 | .02 / .00 | .13 / .02 | .22 | - |
| Situation Awareness & Decision Making | .36 / .23 | .44 / .19 | .09 / .00 | .51 / .23 | .25 | - |
| Mean | .39 / .34 | .32 / .12 | .03 / .00 | .28 / .08 | .19 | - |

765

766  *Note*. In addition to the five-point scale (top), we derived a dichotomous pass/fail scale by assigning

767  the lower two scale-points to *fail* and the upper three scale-points to *pass* (bottom). The mean ICC

768  value was calculated using Fisher z' transformation.

769

770 Table 3. $r_{wg}$ and *ICC(3)* of the LOSA ratings as a function of the performance level shown in the

771 scenario.

772

| LOSA | Performance Level | | | | $r_{wg}$ mean | *ICC(3)* |
|---|---|---|---|---|---|---|
| | Out-standing | Medium-high | Medium-low | Poor | | |
| **Four-point scale** | | | | | | |
| Planning | .76 | .71 | .78 | .68 | .74 | .47 |
| Execution | .78 | .72 | .80 | .74 | .76 | .43 |
| Review / Modify | .70 | .53 | .65 | .65 | .63 | .25 |
| Overall | .70 | .48 | .60 | .65 | .61 | .30 |
| Mean | .74 | .61 | .71 | .68 | .68 | .37 |
| **Pass/fail scale** | | | | | | |
| Planning | .89 | .00 | .59 | .01 | .37 | - |
| Execution | 1 | .03 | .59 | .01 | .41 | - |
| Review / Modify | .89 | .09 | .51 | .07 | .39 | - |
| Overall | .89 | .19 | .69 | .42 | .54 | - |
| Mean | .92 | .07 | .60 | .13 | .43 | - |

773

774 *Note*. In addition to the four-point scale (top), we derived a dichotomous pass/fail scale by assigning

775 the lowest scale point to *fail* and the upper three scale-points to *pass* (bottom). The mean ICC value

776 was calculated using Fisher z' transformation.

777