1     Blind Title Page

2     **Rater agreement of a test battery designed to assess adolescents' resistance training skill**

3     **competency**

4    **Abstract**

5    **Objectives:** The study aim was to assess rater agreement of the Resistance Training Skills Battery

6    (RTSB) for adolescents. The RTSB provides an assessment of resistance training skill competency

7    and includes six exercises. The RTSB can be used to assess performance and progress in adolescent

8    resistance training programs and provide associated feedback to participants. Individual skill scores

9    are based on the number of performance criteria successfully demonstrated and an overall resistance

10   training skill quotient (RTSQ) is created by summing the six skill scores. **Design/Method:** The eight

11   raters had varying experience in movement skill assessment and resistance training and completed a

12   2-3 hour training session in how to assess resistance training performance using the RTSB. The raters

13   then completed an assessment on six skills for 12 adolescents (mean age=15.1 years, SD =1.0, six

14   male and six female) in a randomised order. **Results:** Agreement between seven of the eight raters

15   was high (20 of the 21 pairwise correlations were greater than 0.7 and 13 of the 21 were greater than

16   0.8). Correlations between the eighth rater and each of the other seven raters were generally lower

17   (0.45 to 0.78). Most variation in the assigned RTSB scores (67%) was between cases, a relatively

18   small amount of the variation (10%) was between raters and the remainder (23%) was between

19   periods within raters. The between-raters coefficient of variation was approximately 5%. **Conclusion:**

20   The RTSB can be used reliably by those with experience in movement skill assessment and resistance

21   training to assess the resistance skill of adolescents.

22

23   Word Count: 239

24   2995 for the article

25

26

28

**Introduction**

Youth physical activity guidelines have identified strength as an important health related factor [1], and current public health objectives now aim to increase the number of school-age youth who participate in muscle strengthening activities [2]. Regular participation in an age-appropriate related resistance training program can enhance muscular fitness, power and motor skill performance [3-5]. Furthermore, resistance training interventions in youth can have a positive influence on metabolic health, body composition, cardiorespiratory fitness, blood lipids, bone mineral density and insulin sensitivity [6, 7]. There is clear evidence that resistance training can be a safe, effective and worthwhile activity for children and adolescents provided that appropriate training guidelines are followed and qualified instruction is available [8-10].

Resistance training programs are usually evaluated using 'product' type fitness tests that assess muscular strength and local muscular endurance [11, 12] (i.e. 'how heavy' or 'how many repetitions), rather than providing meaningful feedback on movement skill technique. Movement skill technique is important when assessing the fundamental movement skill competency (i.e. the ability to throw and kick) of children and adolescents as this type of 'process' assessment involves specific feedback regarding which particular components of the skill need to be improved for satisfactory movement skill performance. A process oriented skill assessment involves assessing the 'presence' or 'absence' of a number of components/criteria per skill that are considered essential for mastery of that particular skill. For example, a component of a successful kick is the ability to place the non-kicking foot even with or slightly behind the ball [13].

At present, a process oriented assessment is not commonly used in youth resistance training programs. Therefore, the Resistance Training Skills Battery (RTSB) was developed to assess adolescents' skill competency in resistance training [14]. Potentially, the RTSB could be used to assess each participant's individual performance and, when appropriate, provide general information regarding group level performance and progress in adolescent resistance training programs, while providing constructive feedback to participants. The RTSB includes six skills with each skill involving movements which are considered to provide the basis for strength development. These six skills are summed to provide a resistance training

57    competency total quotient (RTSQ). Initial research was conducted to determine the one week test-retest

58    reliability of the RTSB with 63 adolescents (mean age of 14 years). It was found the RTSB could be used to

59    reliably rank both male and female adolescents on overall resistance training competency and that the RTSB

60    had the necessary sensitivity to detect small changes in resistance skill competency. The RTSB also showed

61    evidence of construct validity,  with the RTSQ predicting 39% of variance in muscular fitness (assessed

62    using handgrip strength, timed push-up and standing long jump tests) [14]. However, the skills in this study

63    were all assessed by the same research assistant, so rater agreement for the RTSB has not been established.

64

65    Rater agreement is the measurement of the consistency or agreement in scores obtained from two or

66    more raters [15, 16], and is important to consider when assessing movement skill proficiency. It is

67    imperative to demonstrate that if a group of raters receive the same training in instrument

68    administration, that they are then able to reliably assess participants' skill competency, otherwise the

69    instrument has limited applicability in the field. Studies of rater agreement in the health literature are

70    often underreported, and when they are reported, they tend to be incomplete and inadequate;

71    therefore, there is a need for such studies to be performed in the future [17].

72

73    When assessing rater agreement it is possible to test the effect of the participant, the rater and also the

74    order of assessment. Analysing for a potential order effect enables an understanding of whether there

75    is a systematic difference occurring during assessment independent from rater differences. For

76    example, if a rater first assesses two adolescents who are poor performers of a skill, the rater as a

77    consequence may then inflate the score of the next adolescent simply because the performance is so

78    much improved from the previous skill performance. Agreement studies that don't test for an order

79    effect are therefore not assessing a potential source of systematic variation. Therefore, the aim of the

80    current study was to assess inter-rater agreement and reliability of the RTSB using the RTSQ.

81    Ordering effects were also assessed.

82

83

**Methods**

Approval for the study was gained from the University Research Ethics Committee and the school principal from one secondary school in New South Wales (NSW), Australia. Parental permission and child assent were obtained. The protocol is described elsewhere [14], but briefly, students completed assessments at school as part of 'all male' or 'all female' groups of three or four. Students observed demonstrations by a research assistant and only questions relevant to the particular exercise (e.g., number of repetitions) were allowed. Encouragement was provided but not skill specific feedback. Students completed two trials of four repetitions for each skill in the following order: (i) body weight squat (ii) push-up (iii) lunge (iv) suspended row (v) standing overhead press and (vi) front support with chest touches. Trunk stability is assessed via *front support* and *chest touches*, upper body pushing strength is assed via a *push-up*, upper body pulling strength is assessed via a *suspended row*, lower body bilateral strength is assessed via a *squat*, and lastly, lower body unilateral strength is assessed via the *lunge*. The exercises therefore target the major muscle groups: lower body (*squat/lunge*), chest, back and arms (*push-up* and *suspended row*), shoulders (*standing overhead press*) and core (*front support with chest touches*). The exercises were all done with only body weight – no additional weight was added. A digital video camera recorded skill attempts. Each skill has four (push-up and suspended row) or five (body weight squat, lunge, standing overhead press and front support with chest touches) performance criteria. Please see Table 1. Scoring was based on the best performance of the skill during the four repetitions for each of the two trials. Participants were awarded a '1' for each criteria correctly demonstrated and '0' if it was not correct. The score for each trial were summed and then totaled for each skill and then the skill scores were all summed for the resistance training skill quotient (RTSQ) (possible range 0 to 56) [14].

TABLE 1 – see supplementary file

For this current study, video assessments of the six skills were selected by taking a stratified random sample of 12 students from the pool of 63 students in the original study (44 males, 19 females, Mean Age 15.1, SD = 1.0). Assessments used for analysis in this manuscript were the first assessments of two trials (assessments were conducted on two occasions seven days apart to determine test retest reliability; this has already been

112   reported [14]). Firstly, all video assessments were grouped by sex and then tertiles were assigned based on the

113   scores assigned previously by the research assistant. Girls and boys performed differently in this original

114   assessment. For girls, the first tertile was a score less than 43 out of the possible 56, the second tertile was

115   from 43 to <47 and the third tertile was $\geq$ 47. For boys, the first tertile was < 40, the second tertile 40 to < 47

116   and the third tertile was $\geq$ 47.  Then two students were randomly selected from each of the six strata.

117

118   Eight raters independently assessed the six videotaped skills for all 12 students (a total of 72 skill

119   assessments per rater). Raters had a range of backgrounds with varying combinations of relevant

120   qualifications, movement skill assessment coding and resistance training experience. Please see Box 1.

121   **Box 1**

| Rater | Relevant Degree/Qualification | Movement skill assessment experience | Resistance training experience |
|---|---|---|---|
| r1 | Physical Education | Extensive experience | 25 years recreational Strength/Conditioning Coach |
| r2 | Physical Education | Limited experience | 10 years recreational |
| r3 | No | Extensive experience | <5 years recreational |
| r4 | Physical Education | Limited experience | <5 years recreational |
| r5 | Physical Education Strength/Conditioning Coach | Limited experience | 10 years recreational |
| r6 | Physical Education | Moderate experience | 8 years recreational |
| r7 | Physical Education | Extensive experience | 8 years recreational |
| r8 | Exercise Science | Little experience | <5 years recreational |

122   Note. Extensive experience = coding >500 performances, Moderate experience = coding >300
123   performances, Limited experience = undergraduate unit, Little experience = a lecture or two.
124

125   Each rater was sent a RTSB training package that included videos for each skill that had been classified in

126   terms of the previous scoring as 'poor' (i.e. few criteria performed correctly), 'medium' (most criteria

127   performed correctly) or 'high' performance (all criteria performed correctly). For example there were three

128   videos of three different students performing the squat to a 'poor', 'medium' or 'high' level. Raters were

129   asked to firstly view these videos and the accompanying scoring sheets which showed how the student had

130   been previously coded. When raters considered they understood the scoring protocol they were asked to

131   code the six skills for each of the 12 students in a specific pre-determined order that was assigned to them.

132    Raters spent on average 90 minutes developing an understanding of the scoring protocol and 120 minutes

133    scoring the trials.

134

135    The order of student assessment (i.e. 1-12 positions) was randomised for each rater. A rater (1 ... 8)

136    was allocated to a presentation order for the assessments by randomly selecting a column from the

137    design matrix for a row-column design.  The row-column design (rows = positions and columns =

138    raters) had the following properties: (1) Each student was assessed once by each rater, (2) Each

139    student was evaluated no more than once in a position, (3) Each pair of students appeared in the same

140    assessment position between 4 and 7 times, and (4) Each student was preceded by every other student

141    no more than once. The design was not balanced for the residual effect (if any) of the evaluation of

142    the preceding student on the evaluation of the current student as this would have required a larger

143    design (such as a Williams' Square) and recruitment of more raters.  Nevertheless the chosen design

144    allowed these residual, or carryover, effects to be estimated.

145

146    As a check on the overall discrimination of the eight raters, a nested analysis of variance (ANOVA)

147    with raters regarded as a random effect and students regarded as a fixed effect, explored whether there

148    was significant variation between the means of the 12 students. In addition, for each student, the

149    variance between the raters was calculated as a check on the stability of the overall assessments and

150    Bartlett's test was used to assess the homogeneity of these within-student (i.e. between-rater)

151    variances. Similarly, for each rater, the variance between the students was calculated as a check on

152    their discrimination and Bartlett's test was used to assess the homogeneity of these within-rater (i.e.

153    between-student) variances. Diagnostic plots of fitted values and residuals were viewed to assess

154    outliers and to check for variance-mean relationships. Agreement between pairs of raters was assessed

155    by computing Pearson's correlation coefficient. The residual effect of the assessment of the previous

156    student on the assessment of a student was investigated via a mixed model analysis (using REML) in

157    which raters (1 to 8) and positions (1 to 12) were regarded as random effects and students, and the

158    previous student (including no previous student, i.e. assessment occurred in the first position), were

159    regarded as fixed effects.  Lastly, in a random effects analysis, variance components for students,

160   raters, and, assessments within raters were estimated to enable intraclass correlations to be reported

161   All analyses were conducted using GenStat Release 14.2 statistical software [18].

162

163   **Results**

164   Mean scores for the 12 cases ranged from 33.9 to 49.75 (Table 2). Table 2 also shows the original

165   tertile assigned to each case (i.e. High/Medium/Low) and the minimum and maximum score assigned

166   for each case by any of the raters. Diagnostic plots of fitted values and residuals showed only one

167   potential outlier (rater 8's relatively low assessment of case #139). The nested ANOVA indicated

168   significant variation between the cases ($p < 0.001$). Two of the 12 cases appeared to have relatively

169   high between-rater variance (or potential discordance), namely cases #130 and #157 and two of the 12

170   cases, namely cases #146 and #139, appeared to have relatively low between-rater variance (or

171   reasonable concordance). The variance also appeared to vary with the mean (lower variances at the

172   high end of the scale where the scores have an upper bound of 56, and higher variances in the middle

173   of the scale, namely 27 to 40). Homogeneity of these between-rater (within-case) variances was

174   explored using Bartlett's test and, despite the apparent differences, there was no significant departure

175   from homogeneity of variance ($\chi^2_{11} = 8.18$; $p = 0.697$).

176

177   TABLE 2

178

179   Agreement between seven of the eight raters was high (20 of the 21 pairwise correlations were greater

180   than 0.7, 13 of the 21 were greater than 0.8 and the range was 0.67 to 0.94). Correlations between the

181   eighth rater (r8) and each of the other seven raters were generally lower (0.45 to 0.78) and this eighth

182   rater also had the highest mean score (Table 3). Mean scores for the eight raters ranged from 37.50 to

183   43.67. Table 3 also shows the maximum and minimum score given by each particular rater.  (Table 3).

184   ANOVA indicated significant variation between the raters ($p < 0.001$).  Two of the raters (r1 and r7)

185   appeared to have relatively high between-case variance, indicating either high discrimination or

186   instability, or, both. One rater (r2) appeared to have relatively low between-student variance,

187   indicating either low discrimination or, moderate to high, stability, or, both.  Homogeneity of these

188    between-student (within-rater) variances was explored using Bartlett's test and, despite the apparent

189    differences, there was no significant departure from homogeneity ($\chi^2_7 = 5.79$; $p = 0.565$).

190

191    The mixed model analysis showed no significant effect of first position (i.e. no previous assessment)

192    versus the other positions ($p = 0.788$) and no overall residual or carryover effect of the assessment of

193    the previous student on the current assessment of a student ($p = 0.411$). When raters (n=8) and

194    students (n=12) were regarded as random effects, the total variance in the 96 RTSB scores was mostly

195    between students (67%), a relatively small amount of the variation (10%) was between raters and the

196    remainder (23%) was between periods within raters (Table 4). The between-raters coefficient of

197    variation was approximately 5%.

198    TABLE 3 and 4

199

200    **Discussion**

201    This study has shown that the RTSB [14] can be used reliably to assess the resistance training skill

202    competency of adolescents. The variation between raters was relatively small, with most of the

203    variation being due to the particular cases that were assessed. Seven of the eight raters commonly had

204    high agreement (pairwise correlations over 0.80). Even the eighth rater (who generally had lower

205    agreement), still had only two pairwise correlations that were below 0.68. Studies which use a process

206    oriented battery to assess the movement skills of children have reported high inter-rater reliability

207    statistics.  For example, a recent Brazilian study involving children reported an ICC of 0.88 for the

208    locomotor subtest and 0.89 for the object control subtest in the Test of Gross Motor Development

209    (TGMD-2) [19]. Similarly, a study of Australian preschool children using the TGMD-2 reported similar

210    results for both subtests (locomotor $ICC = 0.92$ and object control $ICC = 0.90$) [20]. Thus, our estimate

211    of the interrater reliability statistic (ICC = 0.67) for our assessment battery of resistance training skills

212    is lower than such statistics reported in studies of children's movement skill ability that use process

213    oriented instruments.  This could be for several reasons. Firstly, the Brazilian study described their

214    raters as 'expert' and the Australian study reported raters received 12 hours of training, whereas in the

215    current study only three of the eight raters could be called 'expert' (based on a criteria of extensive

216     experience in movement skill assessment combined with some resistance training experience) and the

217     training period was less. Secondly, our study involved eight raters whereas the Brazilian study

218     involved three raters and the Australian study used four raters. Having a higher number of raters

219     purposively selected to have varying levels of experience will increase the observed between-rater

220     variance component and, all other things being equal, decrease the interrater reliability. Furthermore

221     the inclusion of one relatively inexperienced professional whose agreement with the other seven raters

222     was low may have further inflated the between and within rater variance components [21]. In a post-hoc

223     analysis, we excluded the 8[th] rater and found that the ICC measure of rater agreement increased from

224     0.67 to 0.71. Finally it does not appear that either of these studies used a mixed model where potential

225     variance was explained at each potential level (the student/the rater - both between and within) which

226     may also have influenced results.  It has been noted in an article which proposes guidelines for

227     reporting reliability and agreement studies that although ICC values are reported in many health

228     research studies it is often not clear what ICC is being reported and how the analysis has been

229     performed [17]. The same article also suggests that values above 0.60, 0.70, or 0.80 are all reported as

230     minimum values for reliability coefficients, and these values should be seen as appropriate for group-

231     level comparisons and/or research purposes; accordingly, the ICC value found for the current study

232     could be regarded as having met a minimum standard[17].

233

234     Of note, when considering the mean scores for each rater, the raters with less experience coding

235     movement skills had higher overall means than the three raters with considerable experience, even

236     though all raters had relevant backgrounds. The rater with the highest mean score (r8) was the rater

237     with little previous experience. It might be expected that those with experience in observing and

238     coding movement patterns in adolescents would exhibit higher levels of discernment when assessing

239     movement skills and therefore apply more precise scoring. This information may be useful to

240     researchers recruiting movement skill assessors, as well as physical education teachers who may

241     solicit assistance from others during class testing.

242

243 Furthermore 10 of the 12 cases were all rated in the same tertile as those originally assigned, giving

244 further evidence towards the potential of this instrument to be used by a number of raters in a reliable

245 fashion. This study also showed there were no order effects indicating that raters should be able to

246 assess participants in any order and still achieve reliability. However it must be noted that whilst the

247 order of watching and assessment was specified clearly for each rater, the assessment order was not

248 supervised by the researchers. Order effects are not generally reported in literature reporting reliability

249 of movement skill assessment, although one study in preschool children reported that they

250 intentionally ordered the skills for ease of assessment [22]. This study did not however assess any

251 potential order effects in the rater agreement analysis [22].

252

253 **Conclusion**

254 In conclusion, given the high agreement between seven of the eight raters and the relatively low

255 between-rater coefficient of variation, namely 5%, we believe that the RTSB can be used reliably to

256 assess skill competency in selected resistance training exercises in adolescents.

257

258 **Practical Implications**

259 • Raters' with experience in movement skill assessment coupled with at least recreational resistance

260    training experience, are able to reliably assess participants' skill competency after a short training.

261 • The RTSB can be used reliably in adolescent resistance training interventions when supervised by

262    trained assessors with the appropriate backgrounds.

263 • Results from the current study coupled with our previous findings highlight the potential usefulness

264    of the RTSB.

265

270    **References**

271    1       Strong WB, Malina RM, Blimkie CJR, *et al.* Evidence based physical activity for school-age

272    youth. *J Pediatr* 2005; 146: 732-37.

273    2       World Health Organization. *Global Recommendations on Physical Activity for Health.*

274    Geneva 2010; 55.

275    3       Faigenbaum AD, Kraemer WJ, Blimkie CJR, *et al.* Youth Resistance Training:

276    Updated Position Statement Paper from the National Strength and Conditioning Association. *J*

277    *Strength Cond Res.* 2009; 23: S60-79.

278    4       Behringer M, vom Heede A, Matthews M, *et al.* Effects of Strength Training on Motor

279    Performance Skills in Children and Adolescents: A Meta-Analysis. *Pediatr Exerc Sci.* 2011; 23: 186-

280    206.

281    5       Harries S, Lubans DR, Callister R. Resistance training to improve power and sports

282    performance in adolescent athletes: A systematic review and meta-analysis. *J Sci Med Sport.* 2012;

283    15: 532-40.

284    6       Faigenbaum AD, Myer GD. Pediatric Resistance Training: Benefits, Concerns, and Program

285    Design Considerations. *Curr Sports Med Rep.* 2010; 9: 161-68.

286    7       Benson AC, Torode ME, Singh MAF. Effects of resistance training on metabolic fitness in

287    children and adolescents: a systematic review. *Obes Rev.* 2008; 9: 43-66.

288    8       Behm DG, Faigenbaum AD, Falk B, *et al.*  Canadian Society for Exercise Physiology

289    position paper: resistance training in children and adolescents. *Appl Physiol Nutr Metab.* 2008; 33:

290    547-61.

291    9       Lloyd R, Faigenbaum AD, Stone M. Position statement on youth resistance training: the 2014

292    international consensus. British Journal of Sports Medicine 2013(epub ahead of print).

293    10      Lubans DR, Morgan PJ, Aguiar EJ, *et al.* Randomized controlled trial of the Physical Activity

294    Leaders (PALs) program for adolescent boys from disadvantaged secondary schools. *Prev Med.* 2011;

295    52: 239-46.

296    11      Faigenbaum AD, Milliken G, Moulton L, *et al.* Early muscular fitness adaptations in children

297    in response to two different resistance training regimens. *Pediatr Exerc Sci.* 2005; 17: 237-48

298    12    Lubans DR, Sheaman C, Callister R. Exercise adherence and intervention effects of two

299    school-based resistance training programs for adolescents. *Prev Med*. 2010; 50: 56-62.

300    13    Ulrich DA. Test of Gross Motor Development (2nd ed). Austin, TX 2000.

301    14    Lubans DR, Smith J, Harries S, *et al.* Development, Test-Retest Reliability and Construct

302    Validity of the Resistance Training Skills Battery. *J Strength Cond Res*. 9000; Publish Ahead of Print:

303    10.1519/JSC.0b013e31829b5527.

304    15    Goodwin LD. Interrater agreement and reliability. *Meas Phys Educ Exerc Sci*. 2001; 5(1): 13-

305    34.

306    16    Posner KL, Sampson PD, Caplan RA, *et al.* Measuring interrater reliabity among multiple

307    raters: An example of methods for nominal data. *Stat Med*. 1990; 9: 1103-15.

308    17    Kottner J, Audig´e L, Brorson S, *et al.* Guidelines for Reporting Reliability and Agreement

309    Studies (GRRAS) were proposed. *J Clin Epidemiol*. 2011; 64  96-106.

310    18    Payne RW, Harding SA, Murray DA, *et al.* The Guide to GenStat® Release 14 Part 2:

311    Statistics. VSN International: Hemel Hempstead, UK 2011.

312    19    Valentini NC. Validity and Reliability of the TGMD-2 for Brazilian Children. *J Mot Behav*.

313    2012; 44: 275-80.

314    20    Barnett L, Hinkley T, Hesketh K, *et al.* Use of electronic games by young children and

315    fundamental movement skills. *Percept Mot Skills*. 2012; 114: 1023-34.

316    21    Shrout PE, Fleiss JL. Intraclass correlations: Uses in assessing rater reliability. *Psychol Bull*.

317    1979; 86: 420-28.

318    22    Williams HG, Pfeiffer KA, Dowda M, *et al.* A field-based testing protocol for assessing gross

319    motor skills in preschool children: The children's activity and movement in preschool study motor

320    skills protocol. *Meas Phys Educ Exerc Sci*. 2009; 13: 151-65.

321