1  **Reproducibility of objectively measured physical activity: reconsideration needed**

2  Eivind Aadland[1], Ada Kristine Ofrim Nilsen[1], Einar Ylvisåker[1], Kjersti Johannessen[1],

3  Sigmund Alfred Anderssen[1,2]


4  *[1]Western Norway University of Applied Sciences, Faculty of Education, Arts and Sports,*

5  *Department of Sport, Food and Natural Sciences, Campus Sogndal, Norway.*

6  *[2]Norwegian School of Sport Sciences, Department of Sports Medicine, Oslo, Norway.*

7


8


9  **Corresponding author**

10  Eivind Aadland

11  Department of Sport, Food and Natural Sciences, Faculty of Education, Arts and Sports,

12  Western Norway University of Applied Sciences, Campus Sogndal, Box 133, 6851 Sogndal,

13  Norway. Phone: +47 5767 6086; Email: eivind.aadland@hvl.no

14

15

16

17  **Word count main text: 3836; word count abstract: 194**

**Abstract**

Reliability of accelerometer-determined physical activity (PA), and thus the required length of a monitoring period, appears to depend on the analytic approach used for its calculation. We compared reliability of objectively measured PA using different resolution of data in a sample of 221 Norwegian 2-6-year-old children providing 2–3 valid 14-day periods of accelerometer monitoring (ActiGraph GT3X+) during September–October, January–February, and May-June 2015–2016. Reliability (intra-class correlation, ICC) was measured for 1–14 days of monitoring across the measurement periods using linear mixed effect modelling. These results were compared to reliability estimated using different resolution of data using the Spearman Brown formula. The measured reliability improved only marginally with increased monitoring length and levelled-off after 5–6 days. Estimated reliability differed substantially when derived from different resolution of data: 3.9–5.4, 6.7–9.2, 13.4–26.7, and 26.3–87.7 days of monitoring was required to achieve an ICC = 0.80 using an hour-by-hour, a day-by-day, a week-by-week, and a period-by-period approach, respectively. Reliability could not be correctly estimated from any single resolution of data. We conclude that reconsideration is needed with regard to how reproducibility of objectively measured PA is analyzed and interpreted.

**Keywords:** Test-retest; Reliability; Intra-class correlation; Measurement error; Accelerometry

**Introduction**

Procedures used to analyze accelerometry data and criteria applied to define what constitutes a valid physical activity (PA) measurement varies extensively (1). Because behavior vary greatly over time, an important aspect of accelerometer measurements is how many days or periods of measurement that should be included to obtain reproducible estimates of habitual PA levels. Arguably, the "true" habitual PA level would be superior to a short snapshot, as random error in measurements will increase the likelihood of type II errors and thus invalidate study conclusions (2).

Although findings vary between studies in both adults (3-7) and children (8-20), most evidence suggest that a reasonable reliability (i.e., intra-class correlation (ICC)) of ~ 0.70–0.80 are achieved with 3–7 days of monitoring. However, most previous estimates are derived from the Spearman Brown prophecy formula applied to measurements conducted over a single 7-day period. This procedure estimates the number of measurement periods (usually days) needed to obtain a sufficient reliability level, often considered to be an ICC = 0.80, based on variance components and ICC estimates for a single period. Unfortunately, these study designs have received critique for being likely to underestimate the number of monitoring days needed, and their conclusions should therefore be interpreted with caution (21-24). In comparison, studies that have determined the reliability for several periods of measurement over the course of 2 weeks up to a year, have shown considerable intra-individual variation over time (25, 26, 23, 27, 28, 24). Specifically, studies including several seasons have resulted in reliability estimates of ~ 0.50 for one week monitoring in children (26, 23, 24). These findings agree with studies showing substantial seasonal variation in PA in children and adolescents (29-31), which are obviously not captured when relying on a single measurement period.

61       Beyond seasonal variation, there is also differences in reliability between the analytic

62    approaches applied (24, 23). When using a day-by-day approach (estimating reliability from

63    single days of measurement), reliability is estimated from a correction of the residual (within-

64    subject) variance by dividing by the number of scores to be averaged (i.e., the number of

65    monitoring "units" (k), for example days, weeks, etc.) (32). This procedure leads to an

66    underestimation of the residual variance compared to actually measured residual variance

67    over the period (24, 23). Aadland et al (24) determined reliability over a week in a large

68    sample of schoolchildren over 2 seasons, and found a systematic underestimation of residual

69    variance and resulting overestimation of ICCs (0.64–0.77 vs. 0.49–0.63; 14 to 31% difference

70    after controlling for season) using the day-by-day compared to a week-by-week approach.

71    This finding is consistent with findings showing that the reliability of different numbers of

72    monitoring hours per day and days per week is rather similar using a week-by-week approach

73    (27, 24, 28), whereas an increased number of monitoring days inherently will improve

74    reliability when estimated over an increased number of days. Thus, there appears to be a

75    difference in reliability depending on whether the number of measurements needed is

76    estimated from single days and then extrapolated (i.e., using a day-by-day approach) or

77    actually measured over several weeks or periods (i.e., using a week-by-week approach). We

78    infer from this finding that the resolution of data might be fundamentally important for

79    determining reliability. Thus, the resolution of data should be systematically altered, including

80    both higher (i.e., using an hour-by-hour approach) and lower resolution (i.e., using week-by-

81    week and period-by-period approaches) than traditionally applied to thoroughly investigate

82    this hypothesis.

83       The aim of the present study was to extend our previous findings comparing a day-by-

84    day and a week-by-week approach over 2 seasons (24), using a dataset having 2–3 separate

85    14-day periods of monitoring over different seasons in preschool children. Using different

86  resolution of data, we will compare an hour-by-hour, a day-by-day, a week-by-week, and a

87  period-by-period approach to calculate reliability using the same dataset. We hypothesized

88  that reliability for longer periods (up to 14 days) would be overestimated when estimated

89  from higher resolution data (hour-by-hour and day-by-day) compared to using accumulated

90  data over longer measurement periods (week-by-week and period-by-period).

91

92  **Methods**

93  **Participants**

94  The present analysis is based on data obtained in preschool children from the Sogn og

95  Fjordane Preschool Physical Activity Study (PRESPAS) (33), conducted in Norway during

96  2015–2016. Physical activity was measured with accelerometry over one 14-day period in

97  1340 children (September 2015 to June 2016) and over 3 separate 14-day periods in a

98  subsample of 376 children from 3 municipalities (September to October 2015, January to

99  February 2016, and May to June 2016). In the present study, we included all available

100  children for a comparison of "short-term" reliability over 2 consecutive weeks (*cross-*

101  *sectional sample*), and the subsample having repeated measurements for comparison of "long-

102  term" reliability over 2–3 separate periods of measurement (*longitudinal sample*).

103  Our procedures and methods conform to ethical guidelines defined by the World Medical

104  Association's Declaration of Helsinki and its subsequent revisions. The Norwegian Centre for

105  Research Data approved the study protocol. We obtained written informed consent from each

106  child's parents or legal guardian prior to all testing.

107

108  **Procedures**

5

109 Physical activity was measured using the ActiGraph GT3X+ accelerometer (Pensacola, FL,

110 USA) (34). During all measurements, participants were instructed to wear the accelerometer

111 at all times over 14 consecutive days, except during water activities (swimming, showering)

112 or while sleeping (at night). Units were initialized at a sampling rate of 30 Hz. Files were

113 analyzed at 10 second epochs using the KineSoft analytical software version 3.3.80 (KineSoft,

114 Loughborough, UK). Data was restricted to daytime (i.e., hours 06:00 to 23:59). In all

115 analyses, consecutive periods of $\geq 20$ minutes of zero counts were defined as non-wear time

116 (35, 1). Results are reported for overall PA level (cpm), as well as minutes per day spent SED

117 ($< 100$ cpm), in light PA (LPA) (100–2295 cpm), in moderate PA (MPA) (2296–4011 cpm),

118 in vigorous PA (VPA) ($\geq 4012$ cpm), and in moderate-to-vigorous PA (MVPA) ($\geq 2296$

119 cpm), determined using the previously established and validated Evenson et al cut points (36,

120 37). Data were analyzed with wear requirements of $\geq 8$ hours/day and $\geq 3$ weekdays $+ \geq 1$

121 weekend day/week for each separate week. We required 2 valid weeks of measurement for the

122 cross-sectional sample and 4–6 valid weeks for the longitudinal sample ($\geq 2$ periods). As

123 reproducibility is marginally affected by wear hours per day ($\geq 6$ to $\geq 12$ hours/day (27, 24,

124 28), we did not analyze sensitivity to this wear criteria herein.

125

**Statistical analyses**

127 Children's characteristics were reported as frequencies, means and standard deviations (SD).

128 Differences in PA levels between the 3 measurement periods was tested using a mixed effect

129 model including random intercepts for children and including wear time as a covariate.

130 We calculated reliability using 4 approaches based on different resolution of data; 1) hour-by-

131 hour, 2) day-by-day, 3) week-by-week, and 4) period-by-period. Approaches 1, 2 and 3 were

132 applied to the cross-sectional dataset, whereas approaches 1, 2, 3, and 4 was applied to the

longitudinal dataset. Reliability for single hours (hour-by-hour approach), single days (day-by-day approach), single weeks (week-by-week approach), and single periods (period-by-period approach) of measurement ($ICC_s$) were calculated using variance partitioning applying a one-way random effect model not controlling for season (i.e., determining agreement based on an absolute definition) in both samples, whereas a two-way mixed effect model controlling for season (i.e., determining agreement based on a consistency definition) additionally were applied in the longitudinal sample (32). All models were adjusted for wear time by adding wear time as a covariate because wear time has a strong association with PA and SED estimates and also impact reliability (28), and since most studies control for wear time.

We directly determined ("*MEASURED*") reliability for 1–7 monitoring days across 2 consecutive weeks in the cross-sectional dataset (using a week-by-week approach) and for 1–14 monitoring days across 2–3 separate 14-day periods in the longitudinal dataset (using a period-by-period approach). Thus, these analyses is based on the actual variance components for different numbers of monitoring days across weeks and periods. Contrary to this prosedure, we also extrapolated ("*ESTIMATED*") reliability for average measurements ($ICC_k$ = between-subject variance/(between-subject variance + residual variance/$k$)) and the number of measurements needed using the Spearman Brown prophecy formula ($N = ICC_t/(1-ICC_t)*((1-ICC_s)/ICC_s)$, where $ICC_t$ = the desired level of reliability, and $ICC_s$ = the reliability for single measurement) (3, 32). N was rescaled to days ($N_{days}$) for ease of comparison across approaches using mean values of wear hours per day (11.8 hours in cross-sectional dataset; 11.7 hours in the longitudinal dataset), wear days per week (6.3 days in both datasets), and wear days per period (12.6 days in the longitudinal dataset only). The number of measurements (N) needed to obtain a reliable measurement were estimated using an $ICC_t$ = 0.80.

157 In the week-by-week and period-by-period analyses (longitudinal dataset), we additionally

158 calculated 95% limits of agreement (LoA) and coefficients of variation (CV) from the residual

159 variance (i.e., within-subjects) error term based on the variance partitioning models (LoA =

160 $\sqrt{}$residual variance $*\sqrt{2}*1.96$; CV = $\sqrt{}$residual variance/mean values) (38).

161 All analyses were performed using IBM SPSS v. 24 (IBM SPSS Statistics for Windows,

162 Armonk, NY: IBM Corp., USA). A p-value < .05 indicated statistically significant findings.

163

164 **Results**

165 Of the 1340 children included in PRESPAS, 1308 children provided accelerometer data for

166 the cross-sectional analyses, of whom 873 children (52% boys) fulfilled the wear criterion for

167 2 consecutive weeks and were included in the present analysis (Table 1). Of the 376 children

168 included in the longitudinal subsample, 372 provided accelerometer data, of whom 221 (53%

169 boys) had $\geq$ 2 valid measurement periods and were included in the present analysis.

170 The longitudinal analyses included 144 children having 2 measurements and 77 children

171 having 3 measurements across seasons. In general, PA levels were highest during the summer

172 and lowest during the winter (Supplemental Table 1). The greatest differences were seen for

173 VPA (up to 67% difference), overall PA (up to 21% difference), and MVPA (up to 13%

174 difference) (all p < .001), whereas smaller and less consistent differences over seasons were

175 found for other intensities.

176

177 Reliability across 2 consecutive weeks – cross-sectional sample

178    Table 2 shows the reliability of single measurements (ICC) and the ESTIMATED number of

179    monitoring days needed to achieve a reliability of 0.80 (N) using an hour-by-hour, a day-by-

180    day, and a week-by-week approach. The 3 approaches relying on different resolution of data

181    yielded different results; whereas 2.2–4.3 days was needed using an hour-by-hour approach,

182    4.1–7.7 days was needed using a day-by-day approach, and 4.1–14.1 days was needed using a

183    week-by-week approach.

184

185    Table 3 shows the MEASURED reliability over an average of 1 to 7 days of monitoring using

186    a week-by-week approach. Although the pattern of improvement was somewhat different

187    across variables, in general, reliability improved up to a number of 5–6 monitoring days, after

188    which reliability levelled off.

189

190    Reliability across 2–3 separate 14-day periods – longitudinal sample

191    Compared to the results shown for 2 consecutive weeks (Table 2 and 3), the reliability

192    decreased and the required number of monitoring days increased when values were estimated

193    and measured over several seasons (Table 4 and 5). Similar to results based on 2 consecutive

194    weeks, different resolution of data yielded substantially different ESTIMATED values (Table

195    4); whereas 3.9–5.8 days was needed using an hour-by-hour approach, 6.7–10.2 days was

196    needed using a day-by-day approach, 13.4–32.5 days was needed using a week-by-week

197    approach, and 26.3–111.2 days was needed using a period-by-period approach. In contrast to

198    the estimated reliability, MEASURED reliability increased marginally over the first 5–6 days,

199    after which is levelled off (Table 5), similar to the findings in the cross-sectional dataset.

200    Figure 1 shows the estimated (day-by-day, week-by-week, and period-by-period) and

201 measured reliability for 1–14 days of monitoring. The figure shows that the measured

202 reliability is not estimable by the different approaches.

203 Supplemental Figure 1 shows variance components for MVPA. Compared to actually

204 measured variances, the residual variance is underestimated for long monitoring periods by

205 the day-by-day approach and overestimated for short monitoring periods by the period-by-

206 period approach. Increasing the length of the monitoring period also reduced the measured

207 between-subject variance, whereas between-subject variance is kept constant in estimation

208 models.

209 Controlling for season had in general a minor influence on the results, although it influenced

210 reliability for overall PA, VPA and MVPA, for which the seasonal differences were most

211 prominent (Table 2).

212

213 Agreement for 1 week and 1 period of measurement – longitudinal sample

214 Supplemental Table 2 shows 95% LoA and CV for 1 week (week-by-week approach) and 1 period

215 (period-by-period approach) of measurement, indicating to what extent these monitoring periods are

216 capable of capturing PA levels representing one-year habitual activity levels (1 out of 4–6 weeks and 1

217 out of 2–3 periods, respectively). Results were essentially similar for 1 week and 1 period of

218 measurement; CVs were 9–42% across variables, whereas differences up to 332–385 cpm, 91–94

219 minutes/day of SED, 33–37 minutes/day of MVPA, and 17–22 minutes/day of VPA should be

220 expected between monitoring periods over a year.

221

222 **Discussion**

223    The present study aimed to determine and compare the reproducibility of accelerometer-

224    determined PA using different analytic approaches based on different resolution of data over

225    the short-term (2 consecutive weeks in the cross-sectional dataset) and long-term (2–3

226    separate monitoring periods over different seasons in the longitudinal dataset). Our main

227    finding was that reliability of PA as a function of monitoring length, and thus the required

228    number of monitoring days, is not estimable by extrapolation using any single resolution of

229    data. Our findings show that estimation of reliability applying the much-used Spearman

230    Brown formula is invalid, and that reconsideration is needed with respect to the analysis and

231    interpretation of reliability of accelerometry-derived PA measurements.

232    Most previous studies investigating reliability and the required number of accelerometer

233    monitoring days have estimated reliability based on day-by-day analyses using a single 7-day

234    monitoring period (8, 13, 14, 38, 15, 16, 19, 17, 18, 9-12). In general, these studies conclude

235    that 3–7 monitoring days are sufficient in children. In contrast, studies comparing several

236    monitoring periods captured over different seasons, have yielded substantially lower

237    reliability estimates in both adults (25) and children (26, 23, 24), concluding that longer

238    and/or several monitoring periods is needed. Mattocks (26) determined reliability over 4

239    separate 7-day periods over approximately one year using the Actigraph 7164 accelerometer

240    in 11–12-year-old children and found a reliability of 0.45 to 0.59 across variables. Similarly,

241    Wickel & Welk (23) found an ICC of 0.46 over 3 separate 7-day periods to assess steps for

242    the Digiwalker pedometer in 10-year-old children. Finally, Aadland et al (24) found a

243    reliability of 0.29–0.67 across 2 separate periods 3–4 months apart using the Actigraph

244    GT3X+ accelerometer in a large sample of 10-year-old children.

245    The reliability estimates based on a single monitoring period versus several separate periods

246    differ in 2 important ways. Obviously, separate periods are based on measurements collected

247    over a longer time frame, possibly influenced by seasonality, which increase the likelihood of

248 capturing changes in individuals' PA levels over time. These changes over time also cause

249 differences in variance between the monitoring periods, which will attenuate ICCs as the

250 model assumes compound symmetry and the ICC are sensitive to asymmetry (24, 32).

251 Moreover, the statistical analyses are based on different resolution of data; a day-by-day

252 approach for single period data and a week-by-week approach for multiple (weeklong)

253 periods of data. Our results suggest both these differences are influential for the resulting

254 reliability. First, comparable analytic approaches led to lower reliability in the longitudinal

255 dataset than in the cross-sectional dataset (mean ESTIMATED ICC = 0.07 vs. 0.10 for an

256 hour; 0.33 vs. 0.42 for a day, 0.57 vs. 0.77 for a week, respectively; mean MEASURED ICC

257 = 0.51 vs. 0.77 for a week, respectively). These findings show that reliability decreases when

258 more variation is added to the data when capturing a longer time frame with greater variation

259 in behavior. Thus, our findings show that long-term reliability is underestimated when

260 estimated from a single short measurement period. Even more important, our findings suggest

261 that different resolution of data has a major influence on reliability estimates. Adding to the

262 day-by-day (8, 13, 14, 38, 15, 16, 19, 17, 18, 9-12) and the week-by-week approach (26, 23-

263 25) as applied previously, we extended our analysis to include data using higher (hour-by-

264 hour) and lower (period-by-period) resolution, to obtain an even better picture of how data

265 resolution influence reliability. These approaches led to substantially different reliability

266 estimates and numbers of required monitoring days in both samples, particularly in the

267 longitudinal dataset where the number of monitoring days to achieve an ICC = 0.80 based on

268 the hour-by-hour and period-by-period approach varied from (mean) 4.5 to 49.7 days.


269 The differing findings among the analytic approaches based on differing resolution of data

270 result from erroneous estimation of variance components across resolutions. The ICC is

271 calculated from these variance components, which will vary by resolution. The estimated ICC

272 using the Spearman Brown formula will thus be fully dependent on their correct estimation

273  across different resolutions to obtain correct reliability estimates. However, compared to

274  actually measured variances, the residual variance is underestimated for long monitoring

275  periods using high-resolution data and overestimated for short monitoring periods using low-

276  resolution data. Moreover, whereas between-subject variance is kept constant in estimation

277  models, it decreased when stability of data improved over a longer monitoring period. To this

278  end, both variance components underlying the resulting reliability was erroneously estimated

279  compared to those measured when including 1–7 (cross-sectional dataset) and 1–14

280  (longitudinal dataset) monitoring days in a week-by-week and period-by-period analysis,

281  respectively. These results shows that the correct variance components and thus reliability of

282  objectively measured PA as a function of monitoring length is not estimable from any single

283  resolution of data.

284  Previous studies using long-term measurements (i.e., more than a week) have suggested that

285  periods longer than a week and/or several periods are necessary to determine PA reliably (25,

286  27, 24, 28, 26, 23). The findings herein are consistent with these studies in terms of the

287  modest long-term reliability found for a single week (ICC = 0.35–0.64, mean 0.51) and period

288  (ICC = 0.36–0.66, mean 0.52) of measurement. Taken together, our findings and those of

289  others using several separate monitoring periods suggest a typical 3–7-day period of

290  accelerometer monitoring result in a reliability of 0.29–0.67 across variables in children (23,

291  24, 26). Of great importance though, reliability did not improve beyond 5–6 days when

292  measured over 1–14 days. This pattern contrasts reliability estimates derived from the

293  Spearman Brown formula/ICC for average measurements, which are inherently predicted to

294  improve when the number of measurements increase. Thus, our findings indicate a single 7-

295  day measurement protocol would be the best choice in future research, as it maximize

296  reliability and minimize participant and researcher burden. This recommendation is also in

297  line with the results shown for agreement (LoA and CVs), which was similar for a 7-day and

298    a 14-day period. Reliability could possibly be increased by including several separate

299    monitoring periods for each individual, but such an approach would clearly be less feasible

300    for participants as well as researchers.

301    As noise in exposure (x) variables will lead to attenuation of regression coefficients

302    (regression dilution bias), and noise in outcome (y) variables will increase standard errors (2),

303    unreliable measures weaken researchers ability to make valid conclusions in epidemiology.

304    We argue that, in most cases, researchers are interested in the long-term "true" habitual PA

305    level, rather than activity during the most recent days. Although some health characteristics,

306    as for example insulin resistance, lipid metabolism and blood pressure, might change with

307    acute increases or decreases in PA (40), a child's level of fatness, aerobic fitness, or motor

308    skills takes months or years to develop. For such stable traits, association analyses (using PA

309    as an exposure variable) will inherently suffer from regression dilution bias if relying on an

310    insufficient snapshot of children's habitual activity level. For studies evaluating intervention

311    effects (using PA as the outcome variable), low reliability will decrease power. Thus, in both

312    situations, low reliability increase the likelihood of type II errors (2).

313

314    Strengths and limitations

315    The main strength of the present study is the inclusion of a large and representative sample of

316    children and the use of 2 different datasets (cross-sectional and longitudinal) in which 14-day

317    monitoring where used throughout. This allowed for calculation of short- (2 consecutive

318    weeks) and long-term (3 seasons separated by approximately 9 months) reliability using

319    different resolution of data. As reliability estimates depend on the sample variation (41, 38),

320    the validity of the estimated ICCs presented herein should be generalizable to other contexts,

321    including large-scale population studies. Importantly, the use of 14-day monitoring periods

322 allowed for calculation of actual variance components for accumulation of 1–7 and 1–14 days

323 of measurement over the short- and long-term, respectively, and the comparison of these

324 measurements with estimation and extrapolation of these variance components across

325 different periods. Thus, our findings extend those of Aadland et al (24), who directly

326 compared the reliability of children's objectively measured PA using a day-by-day and a

327 week-by-week approach. Importantly, the hour-by-hour approach was included only to test

328 the hypothesis that reliability improved with higher resolution; we find this approach of little

329 practical importance for researchers.

330 Norway has profound seasonal differences in weather conditions and daylight, which may

331 cause changes in PA levels and types across measurement periods. These characteristics

332 might limit generalizability to areas with less pronounced seasonality. Still, as discussed

333 above, our findings are consistent with previous studies when comparing similar approaches

334 for determination of reliability (8, 13, 14, 38, 15, 16, 19, 17, 18, 9-12, 24, 26, 23).

335 Importantly, this seasonal variation will not influence the comparison across the different

336 analytic approaches, as they are based on the same underlying data. Finally, we could have

337 extended our findings by reporting variance partitioning of multiple components (e.g.,

338 participant, day, and season) as shown previously (23), however, such analyses was out of

339 scope for the present paper.

340

341 **Conclusion**

342 We conclude that reliability of objectively measured PA as a function of monitoring length,

343 and thus the required number of monitoring days, is not estimable by extrapolation using any

344 single resolution of data. Our findings suggest the estimation of reliability applying the much-

345 used Spearman Brown formula to a day-by-day approach provide overly optimistic reliability

15

346  estimates and is invalid for estimating reliability over multiple days or periods. Hence, we

347  caution against this practice and recommend future studies measure reliability over separate

348  monitoring periods. Nevertheless, because our results show that reliability levels off after 5–6

349  monitoring days, they support the use of a 7-day measurement protocol. However, the long-

350  term reliability for this protocol in terms of representing the habitual PA level of children

351  across an extended period, is considerably lower than estimated by most previous studies

352  (mean ICC = 0.51–0.52 for 7–14 days of monitoring). These findings strongly indicate

353  reconsideration is needed with respect to the design, analysis, and interpretation of reliability

354  of accelerometry-derived PA measurements.

355

368

369    Competing interests

370    The authors declare that they have no competing interests.

371

372    Data availability

373    The datasets used in the present study are available from the corresponding author on

374    reasonable request.

375

376

377  **References**

378  1. Cain KL, Sallis JF, Conway TL, Van Dyck D, Calhoon L. Using Accelerometers in Youth

379  Physical Activity Studies: A Review of Methods. J Phys Act Health. 2013;10(3):437-50.

380  2. Hutcheon JA, Chiolero A, Hanley JA. Random measurement error and regression dilution

381  bias. BMJ. 2010;340. doi:10.1136/bmj.c2289.

382  3. Trost SG, McIver KL, Pate RR. Conducting accelerometer-based activity assessments in

383  field-based research. Med Sci Sports Exerc. 2005;37(11):S531-S43.

384  doi:10.1249/01.mss.0000185657.86065.98.

385  4. Jerome GJ, Young DR, Laferriere D, Chen CH, Vollmer WM. Reliability of RT3

386  Accelerometers among Overweight and Obese Adults. Med Sci Sports Exerc.

387  2009;41(1):110-4. doi:10.1249/MSS.0b013e3181846cd8.

388  5. Coleman KJ, Epstein LH. Application of generalizability theory to measurement of activity

389  in males who are not regularly active: A preliminary report. Res Q Exerc Sport.

390  1998;69(1):58-63.

391  6. Matthews CE, Ainsworth BE, Thompson RW, Bassett DR. Sources of variance in daily

392  physical activity levels as measured by an accelerometer. Med Sci Sports Exerc.

393  2002;34(8):1376-81.

394  7. Hart TL, Swartz AM, Cashin SE, Strath SJ. How many days of monitoring predict physical

395  activity and sedentary behaviour in older adults? Int J Behav Nutr Phys Act. 2011;8.

396  doi:10.1186/1479-5868-8-62.

397  8. Basterfield L, Adamson AJ, Pearce MS, Reilly JJ. Stability of Habitual Physical Activity

398  and Sedentary Behavior Monitoring by Accelerometry in 6-to 8-Year-Olds. J Phys Act

399  Health. 2011;8(4):543-7.

400   9. Addy CL, Trilk JL, Dowda M, Byun W, Pate RR. Assessing Preschool Children's Physical

401   Activity: How Many Days of Accelerometry Measurement. Pediatr Exerc Sci.

402   2014;26(1):103-9. doi:10.1123/pes.2013-0021.

403   10. Hinkley T, O'Connell E, Okely AD, Crawford D, Hesketh K, Salmon J. Assessing

404   Volume of Accelerometry Data for Reliability in Preschool Children. Med Sci Sports Exerc.

405   2012;44(12):2436-41. doi:10.1249/MSS.0b013e3182661478.

406   11. Hislop J, Law J, Rush R et al. An investigation into the minimum accelerometry wear

407   time for reliable estimates of habitual physical activity and definition of a standard

408   measurement day in pre-school children. Physiol Meas. 2014;35(11):2213-28.

409   doi:10.1088/0967-3334/35/11/2213.

410   12. Penpraze V, Reilly JJ, MacLean CM et al. Monitoring of physical activity in young

411   children: How much is enough? Pediatr Exerc Sci 2006;18(4):483-91.

412   13. Ojiambo R, Cuthill R, Budd H et al. Impact of methodological decisions on accelerometer

413   outcome variables in young children. Int J Obes. 2011;35:S98-S103. doi:10.1038/ijo.2011.40.

414   14. Rich C, Geraci M, Griffiths L, Sera F, Dezateux C, Cortina-Borja M. Quality Control

415   Methods in Accelerometer Data Processing: Defining Minimum Wear Time. Plos One.

416   2013;8(6). doi:10.1371/journal.pone.0067206.

417   15. Kang M, Bassett DR, Barreira TV et al. How Many Days Are Enough? A Study of 365

418   Days of Pedometer Monitoring. Res Q Exerc Sport. 2009;80(3):445-53.

419   16. Murray DM, Catellier DJ, Hannan PJ et al. School-level intraclass correlation for physical

420   activity in adolescent girls. Med Sci Sports Exerc. 2004;36(5):876-82.

421   doi:10.1249/01.mss.0000126806.72453.1c.

422   17. Treuth MS, Sherwood NE, Butte NF et al. Validity and reliability of activity measures in

423   African-American girls for GEMS. Med Sci Sports Exerc. 2003;35(3):532-9.

424   doi:10.1249/01.mss.0000053702.03884.3f.

425  18. Trost SG, Pate RR, Freedson PS, Sallis JF, Taylor WC. Using objective physical activity

426  measures with youth: How many days of monitoring are needed? Med Sci Sports Exerc.

427  2000;32(2):426-31. doi:10.1097/00005768-200002000-00025.

428  19. Janz KF, Witt J, Mahoney LT. The stability of childrens physical-activity as measured by

429  accelerometry and self-report. Med Sci Sports Exerc. 1995;27(9):1326-32.

430  20. Bisson M, Tremblay F, Pronovost E, Julien A, Marc I. Accelerometry to measure physical

431  activity in toddlers: Determination of wear time requirements for a reliable estimate of

432  physical activity. J Sport Sci. 2019;37(3):298-305. doi:10.1080/02640414.2018.1499391.

433  21. Baranowski T, Masse LC, Ragan B, Welk G. How many days was that? We're still not

434  sure, but we're asking the question better! Med Sci Sports Exerc. 2008;40(7):S544-S9.

435  doi:10.1249/MSS.0b013e31817c6651.

436  22. Matthews CE, Hagstromer M, Pober DM, Bowles HR. Best Practices for Using Physical

437  Activity Monitors in Population-Based Research. Med Sci Sports Exerc. 2012;44:S68-S76.

438  doi:10.1249/MSS.0b013e3182399e5b.

439  23. Wickel EE, Welk GJ. Applying Generalizability Theory to Estimate Habitual Activity

440  Levels. Med Sci Sports Exerc. 2010;42(8):1528-34. doi:10.1249/MSS.0b013e3181d107c4.

441  24. Aadland E, Andersen LB, Skrede T, Ekelund U, Anderssen SA, Resaland GK.

442  Reproducibility of objectively measured physical activity and sedentary time over two

443  seasons in children; Comparing a day-by-day and a week-by-week approach. Plos One.

444  2017;12(12). doi:10.1371/journal.pone.0189304.

445  25. Levin S, Jacobs DR, Ainsworth BE, Richardson MT, Leon AS. Intra-individual variation

446  and estimates of usual physical activity. Ann Epidemiol. 1999;9(8):481-8.

447  26. Mattocks C, Leary S, Ness A et al. Intraindividual variation of objectively measured

448  physical activity in children. Med Sci Sports Exerc. 2007;39(4):622-9.

449  doi:10.1249/mss.0b013e318030631b.

450    27. Aadland E, Johannessen K. Agreement of objectively measured physical activity and

451    sedentary time in preschool children. Prev Med Reports. 2015;2:635-9.

452    28. Aadland E, Ylvisåker E. Reliability of Objectively Measured Sedentary Time and

453    Physical Activity in Adults. Plos One. 2015;10(7):1-13. doi:10.1371/journal.pone.0133296.

454    29. Atkin AJ, Sharp SJ, Harrison F, Brage S, Van Sluijs EMF. Seasonal Variation in

455    Children's Physical Activity and Sedentary Time. Med Sci Sports Exerc. 2016;48(3):449-56.

456    doi:10.1249/mss.0000000000000786.

457    30. Gracia-Marco L, Ortega FB, Ruiz JR et al. Seasonal variation in physical activity and

458    sedentary time in different European regions. The HELENA study. J Sports Sci.

459    2013;31(16):1831-40. doi:10.1080/02640414.2013.803595.

460    31. Ridgers ND, Salmon J, Timperio A. Too hot to move? Objectively assessed seasonal

461    changes in Australian children's physical activity. Int J Behav Nutr Phys Act. 2015;12.

462    doi:10.1186/s12966-015-0245-x.

463    32. McGraw KO, Wong SP. Forming inferences about some intraclass correlation

464    coefficients. Psychol Methods. 1996;1(1):30-46. doi:10.1037/1082-989x.1.4.390.

465    33. Nilsen AKO, Anderssen SA, Ylvisåker E, Johannessen K, Aadland E. Physical activity

466    among Norwegian preschoolers varies by sex, age, and season. Scand J Med Sci Sports.

467    2019;29:862-73. doi:10.1111/sms.13405

468    34. John D, Freedson P. ActiGraph and Actical physical activity monitors: a peek under the

469    hood. Med Sci Sports Exerc. 2012;44(1 Suppl 1):S86-S9.

470    35. Esliger DW, Copeland JL, Barnes JD, Tremblay MS. Standardizing and Optimizing the

471    Use of Accelerometer Data for Free-Living Physical Activity Monitoring. J Phys Act Health.

472    2005;2(3):366-83.

473    36. Evenson KR, Catellier DJ, Gill K, Ondrak KS, McMurray RG. Calibration of two

474    objective measures of physical activity for children. J Sports Sci. 2008;26(14):1557-65.

475    doi:10.1080/02640410802334196.

476    37. Trost SG, Loprinzi PD, Moore R, Pfeiffer KA. Comparison of Accelerometer Cut Points

477    for Predicting Activity Intensity in Youth. Med Sci Sports Exerc. 2011;43(7):1360-8.

478    doi:10.1249/MSS.0b013e318206476e.

479    38. Weir JP. Quantifying test-retest reliability using the intraclass correlation coefficient and

480    the SEM. J Strength Cond Res. 2005;19(1):231-40.

481    39. Chinapaw MJM, de Niet M, Verloigne M, De Bourdeaudhuij I, Brug J, Altenburg TM.

482    From Sedentary Time to Sedentary Patterns: Accelerometer Data Reduction Decisions in

483    Youth. Plos One. 2014;9(11). doi:10.1371/journal.pone.0111205.

484    40. Thompson PD, Crouse SF, Goodpaster B, Kelley D, Moyna N, Pescatello L. The acute

485    versus the chronic response to exercise. Med Sci Sports Exerc. 2001;33(6 Suppl):S438.

486    41. Bland JM, Altman DG. Statistical methods for assessing agreement between two methods

487    of clinical measurement. Lancet. 1986;1(8476):307-10.

488

**Figure legends**

**Figure 1. Measured and estimated reliability for MVPA over 1–14 monitoring days across 3 seasons.** The measured reliability is calculated using a period-by-period approach by accumulating and averaging MVPA over 1–14 monitoring days for each period, thus, the model is based on actual variances. The estimated reliability is calculated for 1 day (day-by-day approach), 1 week (week-by-week approach), and 1 period (period-by-period approach) and extrapolated over k days. All results are based on reliability estimates for a two-way mixed model controlling for season (i.e., a consistency definition of reliability) in addition to wear time.

**Supplemental Figure 1**. **Measured and estimated variance components for MVPA over 1–14 monitoring days across 3 seasons**. The between-subject variance is the part of the variance explained by subjects ("true" variation), whereas the residual variance is the unexplained variance (within-subjects variance or error).

**Supplemental Table 1**. Physical activity levels over 3 seasons (longitudinal sample, n = 221).

**Supplemental Table 2**. 95% limits of agreement and coefficients of variation for 1 week and 1 period of measurement.