

Cross-validation and comparison of energy expenditure prediction models using count-based and raw accelerometer data in youth

Abstract

Machine learning may improve energy expenditure (EE) prediction from body-worn accelerometers. However, machine learning models are rarely cross-validated in an independent sample, and the use of machine learning raises additional questions including the effect of accelerometer placement and data type (count vs. raw) for optimal EE prediction. **Purpose:** To assess the accuracy of artificial neural network (ANN) models for EE prediction in youth using count-based or raw data from accelerometers worn on the hip, wrist, or in combination and compare these to count-based, EE regression equations. **Methods:** Data were collected in two settings, one (n=27) to calibrate the EE prediction models, and the other for model cross-validation (n=34). Participants wore a portable metabolic analyzer (EE criterion) and accelerometers on the left wrist and right hip while completing 30 minutes of exergames (calibration, cross-validation) and a maximal exercise test (calibration only). Six ANNs were created from the calibration data, separately by accelerometer placement (hip, wrist, combination) and data format (count-based, raw) to predict EE (15-second epochs). Three count-based linear regression equations were also developed for comparison to the ANNs. **Results:** The count-based, hip ANN demonstrated lower error (RMSE: 1.2 METs) than all other ANNs (RMSE: 1.7-3.6 METs) and EE regression equations (RMSE: 1.5-3.2 METs). However, all models showed bias toward the mean. **Conclusion:** An ANN developed for hip-worn accelerometers had higher accuracy for EE prediction during an exergame session than wrist or combination ANNs, and ANNs developed using count-based data had higher accuracy than ANNs developed using raw data.

Keywords: Activity trackers, physical activity, machine learning, pattern recognition, out-of-sample

1 **Introduction**

2 Despite the well-known benefits of physical activity (PA) participation in youth,
3 the majority do not meet recommended PA levels (Esteban-Cornejo, Tejero-Gonzalez,
4 Sallis, & Veiga, 2015; Troiano et al., 2008; US Department of Health and Human
5 Services, 2018). Measurement of energy expenditure (EE) using accelerometers is
6 common for determining the volume and intensity of PA, and accurate EE measurement
7 is critical for identification of, and intervention in, youth with low PA. Due to memory
8 capacity and battery life limitations, early accelerometers summarized raw data into
9 ‘activity counts’ or other condensed storage forms on-board the accelerometer in 1-60+
10 second intervals (epochs), meaning that raw data were not available for download. Newer
11 accelerometers allow access to raw (*g*) data collected at high sampling rates for days or
12 weeks at a time (John & Freedson, 2012).

13 EE prediction models developed for count-based accelerometer data are
14 inherently limited in their applicability to other accelerometer brands because counts are
15 brand-specific and often proprietary (John & Freedson, 2012). The use of raw data has
16 the potential to improve the application of models across accelerometer brands and give
17 transparency to features and models used to interpret accelerometer data and associated
18 outcomes (van Hees et al., 2014; van Hees et al., 2013). However, the majority of studies
19 that have developed models for EE prediction have relied on raw data metrics (e.g., mean
20 and percentiles of signal) which are subject to orientation-dependency and are, therefore,
21 influenced by factors such as the angle, attachment method, and side of the body on
22 which an accelerometer is worn. Conversely, count-based data are non-negative and
23 cumulative and, thus, less likely to be influenced by device orientation. Such differences

24 in these types of data may make predictive models using orientation-dependent raw data
25 more prone to over-fitting (Montoye, Pivarnik, Mudd, Biswas, & Pfeiffer, 2016). The
26 vector magnitude (VM) of triaxial accelerometers has been proposed as a strategy to
27 alleviate such issues of orientation dependency for both count and raw data and has been
28 used with hip-worn accelerometers and accelerometers placed on alternate locations
29 (Sasaki, John, & Freedson, 2011; van Hees et al., 2013). However, whether VM improves
30 EE prediction accuracy compared to using triaxial data is equivocal, especially when data
31 type (count vs. raw) and accelerometer placement are considered (Montoye et al., 2016).

32 Several types of predictive models have been developed for translating
33 accelerometer data into EE in youth, ranging in complexity from count-based regression
34 models (Crouter, Horton, & Bassett, 2012; Freedson, Pober, & Janz, 2005) to machine
35 learning models (Mackintosh, Montoye, Pfeiffer, & McNarry, 2016; Trost, Wong,
36 Pfeiffer, & Zheng, 2012), which use count-based or raw data as inputs. Machine learning
37 models have generally yielded more accurate predictions of EE than linear regression
38 models in initial calibration settings in both youth (Mackintosh et al., 2016; Trost et al.,
39 2012) and adult samples (Montoye, Begum, Henning, & Pfeiffer, 2017). However,
40 several studies in adults (Gyllensten & Bonomi, 2011; Lyden et al., 2014; Sasaki et al.,
41 2016; Staudenmayer et al., 2015) and one study in youth (Hibbing, Ellingson, et al.,
42 2018) have demonstrated that the accuracy of these models decreases when cross-
43 validating in a new or independent sample, indicating a tendency for machine learning
44 models to be over-fit to the data. Further research is, therefore, required to determine the
45 accuracy of machine learning models for the prediction of EE in independent data sets.

46 Finally, it must be considered that the potential benefit of raw data and/or
47 machine learning modeling for EE prediction may be dependent on number and
48 placement of accelerometers used. While the hip is the most common accelerometer
49 placement, wrist-worn accelerometers have seen increased use in recent years due to
50 improved wear-time compliance and ability to capture behaviours such as activity type
51 and sleep (Montoye, Moore, Bowles, Korycinski, & Pfeiffer, 2016; Troiano, McClain,
52 Brychta, & Chen, 2014). Indeed, early research utilizing count-based data has shown
53 poorer accuracy of wrist-worn, compared to hip-worn, accelerometers in adults (Bouten,
54 Sauren, Verduin, & Janssen, 1997; Swartz et al., 2000). More recent studies in youth and
55 adult samples indicate that EE prediction from wrist-worn accelerometers is improved
56 when more complex modeling approaches and/or triaxial/VM data were used instead of
57 vertical axis data and/or simple linear regression models on activity count data (Crouter,
58 Flynn, & Bassett, 2015); Montoye et al., 2017; O'Driscoll et al., 2018). Conversely, EE
59 prediction from hip-worn accelerometers may be less affected by modeling method
60 (Montoye et al., 2017). Additionally, use of multiple accelerometers placed on different
61 body locations sometimes, but not always, leads to improved prediction accuracy
62 (Mackintosh et al., 2016). More research is needed to understand the accuracy of
63 accelerometers worn on various body locations.

64 Given the current gaps in our understanding of whether data type (count vs. raw),
65 accelerometer number (one vs. multiple), and modeling method (machine learning vs.
66 linear regression) affect EE prediction accuracy in youth, the present study's aims were to
67 use an independent sample, cross-validation design to 1) determine if count-based or raw
68 data inputs into a machine learning model yield better EE prediction accuracy, 2)

69 investigate whether hip- or wrist-worn accelerometer data as model inputs (or a
70 combination thereof) yield higher EE prediction accuracy, and 3) compare the accuracy
71 of these machine learning models to three count-based EE prediction regression
72 equations.

73 **Methods**

74 In the present study, we describe the development (calibration) of six artificial
75 neural networks (ANNs) and then focus on the EE prediction accuracy of these ANNs
76 and three count-based regression models in an independent, cross-validation setting. Each
77 institution's respective ethics board approved this study, and participants and
78 parents/guardians provided assent and consent prior to completing the study, respectively.

79 *Calibration participants and protocol*

80 Descriptions of the sample and procedures used for the calibration portion of this
81 study have been described previously (Mackintosh et al., 2016). Briefly, 27 youth (15
82 boys; 11.6 ± 1.0 years) from Swansea, UK participated in a protocol in which they played
83 active video games (exergames; two sessions of 15 minutes with a break between
84 sessions; games included River Rush and Kinect Adventures Reflex Ridge on Xbox 360)
85 and an incremental, graded treadmill test to volitional exhaustion. During the protocol,
86 participants wore an ActiGraph wGT3X-BT (ActiGraph Corp., Pensacola, FL, USA)
87 accelerometer at the right hip and left wrist (collecting raw, triaxial data at 100 Hz) and a
88 METAMAX 3B (Biophysik, Leipzig, Germany) metabolic analyzer (collecting breath-
89 by-breath oxygen consumption).

90 *Cross-validation participants and protocol*

91 Participants in the cross-validation study were 34 youth (Table 1; 21 boys; 10.3 ±
92 1.1 years) from the community of East Lansing, MI, USA. Participants were free from
93 any metabolic or physical condition that would alter their ability to perform, or alter their
94 metabolic response to, the study protocol.

95 ****TABLE 1 HERE****

96 Participants were asked to visit the laboratory for a single visit at least two hours
97 postprandial and having avoided caffeine and strenuous exercise for at least 24 hours
98 prior to their visit. Initially, stature was measured to the nearest 0.01 m (Harpenden
99 stadiometer, Holtain, Crymych, United Kingdom) and body mass to the nearest 0.1 kg
100 (Seca digital scale, Hamburg, Germany) using standardized procedures.

101 Subsequently, two exergames were selected at random from a list of four games
102 previously shown to elicit moderate- to vigorous-intensity PA [MVPA; Kinect
103 Adventures Reflex Ridge, Just Dance 3, Wipeout, and Kinect Sports Boxing; (Barkman,
104 Pfeiffer, Diltz, & Peng, 2016; Clevenger & Howe, 2015; Rosenberg et al., 2013)]. The
105 two games were completed on the easiest level in both single- and multi-player mode
106 (with a research assistant or friend/sibling) for 15 minutes each, resulting in four
107 conditions (two games x two modes). Between games, participants were provided with a
108 5-10 minute break (e.g., to drink water, use the restroom).

109 Throughout all sessions and breaks (with the exception of participants consuming
110 water or using the restroom), breath-by-breath gas exchange was assessed (Oxycon
111 Mobile, Carefusion, Yorba Linda, CA, USA). The Oxycon has been shown to provide

112 reliable and accurate measures of oxygen consumption compared to the Douglas bag
113 method (Rosdahl, Gullstrand, Salier-Eriksson, Johansson, & Schantz, 2010).
114 Additionally, two ActiGraph GT3X+ accelerometers were worn, one on the right hip at
115 the level of the anterior axillary line (orientation: y-axis vertical, x-axis medial-lateral, z-
116 axis anterior-posterior) and one on the posterior aspect of their left wrist between the
117 styloid processes of the radius and ulna (orientation: y-axis vertical, x-axis medial-lateral,
118 z-axis anterior-posterior when in anatomical position); both were secured in place with
119 elastic belts. The ActiGraph GT3X+ has a range of ± 6 gravitational (g) units, and all
120 monitors were set to collect triaxial data in raw mode at a sampling rate of 30 Hz.

121 *Data processing*

122 Data from the exergames, rest intervals, and for calibration, the incremental
123 (graded) treadmill test were used for this study. Thus, the calibration and cross-validation
124 protocols included both steady-state and non-steady-state data.

125 For the calibration data only, the 100 Hz accelerometer data were downloaded in
126 raw form and downsampled to 30 Hz to avoid issues in data comparability between
127 calibration and cross-validation given previous work showing that the use of different
128 sampling rates affects the conversion of raw data to activity counts (Brønd & Arvidsson,
129 2016). To downsample the data, the downloaded .gt3x files from the accelerometer were
130 converted to .wav files using in-house Java software (Oracle Corp., Redwood Shores,
131 CA). These files were subsequently read into MATLAB (MathWorks Inc., Natwick, MA)
132 and resampled to 30 Hz using the *resample* function available in MATLAB (Lyons,

133 2013). Once resampled, the 30 Hz files were converted back to .gt3x files using the same
134 Java program. All subsequent analyses were conducted with 30 Hz data.

135 For both calibration and cross-validation, six features (mean and variance from
136 each of the three accelerometer axes) were calculated from the raw accelerometer data in
137 15-second epochs using the feature extraction tool in ActiLife version 6.13 software
138 (ActiGraph Corp., Pensacola, FL, USA). Data were also downloaded as activity counts in
139 1-second epochs using the ActiLife software, and six features (mean and variance of the
140 activity counts in each accelerometer axis) were calculated from this count data in 15-
141 second epochs using Microsoft Excel 2013 (Microsoft Inc., Redmond, WA, USA). These
142 features were chosen in accord with previous research developing machine learning
143 models to predict EE in youth (Mackintosh et al., 2016), and 15-second epochs were
144 chosen as previous research has shown the transient activity patterns of youth may
145 necessitate shorter epochs than the traditional 60-second epochs used in adults (Bailey et
146 al., 1995).

147 Relative oxygen uptake data from the metabolic analyzers were downloaded in
148 15-second epochs and converted to corrected metabolic equivalents (METs). Specifically,
149 equations adapted from Schofield (1985) were used to predict basal metabolic rate in
150 kcals/day. Next basal metabolic rate was converted to milliliters of oxygen consumed per
151 minute and subsequently to ml/kg/min for determination of age and sex-specific youth
152 metabolic equivalent (MET) values (FAO/WHO/UNU, 2001; Schofield, 1985). For
153 example, if a participant's basal metabolic rate was predicted to be 4.0 ml/kg/min, then
154 Schofield-corrected METs were calculated for this participant by dividing their relative
155 oxygen consumption data by 4.0 (rather than 3.5, which is common in adults). This

156 procedure is supported by a position statement published by the CDC/NIC/NCCOR
157 Research Group on Energy Expenditure in Children (McMurray et al., 2015).

158 Once accelerometer data features and criterion corrected METs were calculated,
159 they were time-aligned. All 15-second epochs where criterion EE was <0.5 corrected
160 METs were removed as this generally represents non-wear or poor sampling (e.g.,
161 occluded sample line) in a given epoch (Mackintosh et al., 2016); this resulted in removal
162 of $\sim 1.4\%$ of the epochs in the calibration dataset and $\sim 0.3\%$ of the epochs in the cross-
163 validation dataset.

164 Using the features calculated from the count-based and raw data, six ANNs were
165 created (using calibration data) and then tested (using cross-validation data) to predict
166 EE; these ANNs were hip count, wrist count, combination count, hip raw, wrist raw, and
167 combination raw, wherein “combination” used a combination of both hip and wrist data.
168 The ANNs included in this study were feedforward, had one hidden layer with five
169 hidden units, and did not have skip-layer connections; all ANNs were developed using
170 the *nnet* package in the R software (Ripley & Venables, 2016). Access to sample data and
171 the ANNs can be found at the following link:
172 <https://drive.google.com/open?id=1SlnXJBh6WUpxJJAjAovVbNw8hW54PhbZ>. ANNs
173 were chosen instead of other machine learning models since previous research shows
174 promise for their use in EE estimation from accelerometer data (Mackintosh et al., 2016;
175 Montoye, Mudd, Biswas, & Pfeiffer, 2015; Preece et al., 2009; Staudenmayer, Pober,
176 Crouter, Bassett, & Freedson, 2009).

177 Additionally, our study sought to compare our developed ANNs to traditional,
 178 regression-based EE prediction methods. However, no previous work has developed
 179 regression equations to predict EE from count-based, hip- or wrist-worn ActiGraph
 180 accelerometer data in 15-second epochs, and there are indications that epoch length
 181 affects accelerometer output (McClain, Abraham, Brusseau, & Tudor-Locke, 2008).
 182 Therefore, we developed (using calibration data) and tested (using cross-validation data)
 183 three in-house regression equations for predicting EE from the VM activity counts
 184 according to accelerometer placement (hip, wrist, or combination). These equations were
 185 developed in SPSS version 24.0 (IBM Corp., Armonk, NY, USA). The resulting
 186 equations are as follows, where “HVM” signifies VM counts ($VM = \sqrt{x^2 + y^2 + z^2}$)
 187 from the hip accelerometer per 15 seconds and “WVM” signifies VM counts from the
 188 wrist accelerometer per 15 seconds:

- 189 • E1: Hip: METs = 0.002346*HVM + 2.576510
- 190 • E2: Wrist: METs = 0.000898*WVM + 2.495456
- 191 • E3: Hip and wrist combination: METs = 0.001078*HVM + 0.000591*WVM +
 192 2.339118

193 *Cross-validation data analysis*

194 All data violated tests for normality, so non-parametric statistics were used. For
 195 each of the nine modeling approaches (six ANNs and three regression equations),
 196 predicted EE was averaged across epochs, separately for each activity (rest/transition and
 197 in each of the four exergames). Predicted EE from each model for each activity was
 198 compared to the criterion using a related-samples Friedman analysis of variance. In the

199 event of a significant overall test statistic, *post hoc* differences between model predictions
200 and the criterion were evaluated using pairwise, related-samples Wilcoxon rank sum
201 tests.

202 For each epoch, squared error was calculated to compare each of the nine
203 modeling approaches (six ANNs, three regression equations) to the criterion EE. Then,
204 for each participant, root mean squared error (RMSE) was calculated, separately for each
205 model. A related-samples Friedman analysis of variance test was used to compare RMSE
206 across ANN and regression models, with *post hoc* differences between models evaluated
207 using a pairwise, related-samples Wilcoxon rank sum test. A p-value <0.05 was used to
208 indicate statistical significance, and a false discovery rate correction was used to account
209 for multiple comparisons (Glickman, Rao, & Schultz, 2014). Bland-Altman plots (Bland
210 & Altman, 1986) were also created to evaluate bias in EE prediction across the nine
211 models evaluated. These plots revealed an outlier for the three raw data ANNs. To
212 determine if the outlier data affected our findings, we ran both Friedman analyses (the
213 analysis comparing criterion EE to predicted EE for each activity and the analysis
214 comparing RMSE among model types) two times, once with outlier data included and
215 once with outlier data excluded. There were no changes in the statistical significance of
216 the findings of either Friedman analysis, so data in the Results are shown with the outlier
217 data included. Analyses were conducted using SPSS version 24.0.

218 **Results**

219 RMSE for EE prediction for each modeling approach is shown in Figure 1. The
220 Friedman test statistic was statistically significant; *post hoc* analyses for RMSE revealed

221 that the count-based, hip ANN had significantly lower RMSE than all other ANNs and
222 regression models. More specifically, the count-based hip ANN had RMSE (mean \pm
223 standard deviation of 1.2 ± 0.3 METs) 19.2% lower than the next best model (hip
224 regression; 1.5 ± 0.5 METs). Conversely, the raw wrist ANN (RMSE: 3.6 ± 1.7 METs) had
225 significantly higher RMSE than all other ANNs and all regression models except the
226 wrist regression model (RMSE: 3.2 ± 0.8 METs). For both the hip and wrist ANNs, the
227 count-based models had significantly lower RMSE (39.6-41.0% lower) than the raw
228 models, although this was not the case for the combination ANNs. In comparing ANNs to
229 the regression models, the count-based ANNs for each accelerometer placement (and
230 combination) had significantly lower RMSE than their corresponding regression
231 equations, and the raw ANNs had significantly higher RMSE than their corresponding
232 regression equations.

233 ****FIGURE 1 HERE****

234 Bland-Altman plots (Figure 2) revealed “bias toward the mean”, where all nine
235 predictive models overestimated EE when criterion-measured EE was low (i.e., during
236 low-intensity activities) and underestimated EE when criterion-measured EE was high
237 (i.e., during high-intensity activities). Additionally, these plots revealed narrower 95%
238 limits of agreement for all three count-based ANNs compared to the raw ANNs for the
239 hip, wrist, and combination ANNs. Limits of agreement for the hip, wrist, and
240 combination regression models were wider than that of the count-based ANNs but
241 narrower than the raw data ANNs. As indicated in the Methods, these plots revealed an
242 outlier in the dataset for one participant, for which EE was substantially overestimated by
243 all three raw ANNs (average overestimation of 4.6 METs for hip, 10.3 METs for wrist,

244 and 6.0 METs for combination). This participant's data were removed and the data
245 reanalyzed, which reduced RMSE to 1.9 METs (from 2.0 METs) for the raw hip ANN,
246 3.4 METs (from 3.6 METs) for the raw wrist ANN, and 1.7 METs (from 1.9 METs) for
247 the raw combination ANN, but there was no change in the overall findings. Similarly,
248 removal of the outlier lowered the limits of agreement (shown as low, high) for the raw
249 hip, wrist, and combination ANNs [(-4.3, 2.2 METs), (-4.5, 7.8 METs), and (-3.8, 2.3
250 METs), respectively] but did not affect overall comparisons.

251 ****FIGURE 2 HERE****

252 Criterion-measured and accelerometer-predicted EE for each activity in cross-
253 validation are shown in Table 2. Criterion data from the calibration dataset were
254 comparable to that of cross-validation, with an EE of 2.8 ± 1.6 METs during
255 rest/transitions, 3.7 ± 1.2 METs during the exergame Kinect River Rush, 4.6 ± 1.6 METs
256 during the exergame Kinect Adventures Reflex Ridge, and 6.5 ± 2.6 METs during the
257 treadmill test in calibration. In line with the overall RMSE analysis, the count-based hip
258 ANN performed best; while predicted EE was significantly different from the criterion at
259 rest and for two (of four) exergames, average EE predictions were within 0.6 METs of
260 the criterion measure for all activities and were not different from the criterion overall.
261 Average EE from the hip regression equation was the next best, with EE within 0.7
262 METs of the criterion for all activities (although all were statistically different from the
263 criterion). Conversely, both wrist ANNs, the count-based combination ANN, the wrist
264 regression equation, and the hip-wrist combination regression equation significantly
265 overestimated EE overall and for all activities, with biases of 1.4 to 2.2 METs overall, 0.5
266 to 1.6 METs for rest/transitions, and 0.4 to 4.1 METs during the exergames.

267 **TABLE 2 HERE**

268 **Discussion**

269 Our study used a semi-structured setting primarily involving exergame play to
270 determine whether data type (count-based or raw) and/or accelerometer placement (hip,
271 wrist, or combination) affect ANN-based EE prediction accuracy in youth and how such
272 EE prediction models compared to count-based regression equations. Overall, an ANN
273 machine learning model using count-based, hip accelerometer data had lower error in
274 predicting EE compared to ANNs developed using wrist or combination data, with
275 RMSE 29.2% lower than the next best performing ANN. Notably, count-based ANNs
276 generally outperformed raw data ANNs. Additionally, the count-based, hip ANN had
277 lower error than three count-based EE regression equations (MAPE 19.2% lower than the
278 best performing regression model), which is consistent with past work showing that
279 machine learning methods may improve EE prediction compared to simple regression
280 models (Montoye et al., 2015; Staudenmayer et al., 2009).

281 We report lower error of predictive models using count-based, hip-worn
282 accelerometer data compared to wrist- or combination of hip- and wrist-worn
283 accelerometers. Comparisons of models using data from hip- versus wrist-worn
284 accelerometers for predicting EE have had mixed results, with studies across youth and
285 adult samples indicating lower RMSE values from either hip (Hibbing, LaMunion,
286 Kaplan, & Crouter, 2018; Mackintosh et al., 2016) or wrist (Crouter et al., 2015; Ellis et
287 al., 2014; Staudenmayer et al., 2015) wear locations, or indicating that which model
288 performed better depended on the input features (Montoye et al., 2015). Of note, previous

289 studies have generally reported small differences in RMSE between hip-and wrist-worn
290 monitors (e.g., 0.1-0.2 METs), in contrast to the larger differences in RMSE in the
291 present study (1.2-3.6 METs).

292 The larger RMSE values in the present study compared to previous studies may
293 be due to our focus on exergames, which involve sporadic arm movements. Graves et al.
294 (2008) previously found that hip-worn accelerometers could better predict EE than
295 accelerometers placed on the upper limb during exergames, indicating that the higher
296 accuracy of hip-worn models may be due to the types of activities included in the present
297 study. Additionally, Hwang, Fernandez, & Lu (2018) reported poorer reliability for
298 ActiGraph monitors worn on the wrist compared to the hip during exergames, further
299 supporting that the poorer performance of the wrist-worn monitor may be at least
300 partially due to the focus of the present study on exergames.

301 It also may be that wrist choice matters when wearing an accelerometer during
302 exergame play. Graves et al. (2008) found that non-dominant arm movement was largely
303 impacted by the type of exergame youth participate in and/or their skill level, while
304 dominant arm movement was not. While we did not assess participant handedness in our
305 samples, population estimates suggest that only ~8% of individuals are left-hand
306 dominant (McManus, 1991), so our left-wrist accelerometer was likely the non-dominant
307 wrist for the vast majority of the sample. Future research should therefore ascertain
308 whether an accelerometer worn on the dominant wrist would be preferable for improving
309 EE prediction accuracy in this setting.

310 Better accuracy of the hip ANN in the independent sample in the present study
311 may also be due to less movement variability at the hip compared to the wrist among
312 participants and among different exergames, which is likely also the case with other non-
313 ambulatory activities that take place in free-living settings. Variability in wrist
314 movements during this type of activity may have also contributed to the poorer accuracy
315 of the wrist and combination in this independent sample cross-validation compared to our
316 previous calibration study in the same setting and the same population (Mackintosh et al.,
317 2016), whereas the hip ANNs were affected but to a lesser degree. Despite interest in
318 wrist-worn accelerometers for the purposes of improved compliance (Troiano et al.,
319 2014) and/or measurement of other health-related behaviours such as sleep (van Hees et
320 al., 2015), more work is needed to improve EE prediction accuracy of wrist-worn activity
321 monitors. A recent meta-analysis by O’Driscoll et al. (2018) suggests that the accuracy of
322 wrist- and arm-worn monitors for predicting EE is improved with addition of
323 physiological data such as heart rate, so future work should evaluate this and other
324 additional sensing methods as a potential way to improve wrist-based EE prediction.

325 A second important finding is that ANNs developed independently from hip- and
326 wrist-worn, count-based data had lower error than corresponding raw data ANNs. This
327 may be related to count-based data being designed specifically to capture acceleration
328 frequency/magnitude and to filter accelerations that occur outside of a certain range
329 (ActiGraph, 2016; Brønd & Arvidsson, 2016). Indeed, the conversion of raw data into
330 counts may reduce instances where aberrant movements unduly affect EE prediction.
331 However, despite the superior performance of count-based ANNs compared to raw data
332 ANNs in the present study, it is pertinent to note that counts are a manufacturer-specific

333 metric that cannot easily be translated or compared across accelerometer brands (John &
334 Freedson, 2012), contrary to raw data which should be similar. The proprietary nature of
335 count generation and the non-comparability of count-based data across brands render
336 count-based models of limited use, unless 1) ActiGraph monitors are used or 2) count
337 data equivalent to the ActiGraph are generated from the raw data of other accelerometer
338 brands, which is now possible due to recent work by Brønd et al. (2017). Nonetheless, the
339 higher accuracy of count-based than raw data models found in this study is informative
340 and may offer researchers information as to how to improve the accuracy of raw data
341 modeling techniques. Future research could investigate filtering methods, other features
342 such as frequency-domain features, and possible translation of raw data into orientation-
343 independent metrics such as VM, Euclidean norm minus one, or mean amplitude
344 deviation. Such methods may allow for the use of the meaningful aspects of raw data
345 while also removing signal noise (Bai et al., 2016; Bakrania et al., 2016). Additionally,
346 by making these methods open-access, the comparability of data across brands would be
347 preserved, allowing predictive models to be used across accelerometer brands.

348 A final notable finding is that the ANNs developed from a combination of hip and
349 wrist data had poorer accuracy than the hip ANNs and hip regression equations in our
350 study. Findings comparing accuracy of single- and multiple-accelerometer prediction
351 methods are mixed. For example, two studies by Dong et al. (Dong, Biswas, Montoye, &
352 Pfeiffer, 2013; Dong, Montoye, Moore, Pfeiffer, & Biswas, 2013) found that a three-
353 accelerometer system (wrist, thigh, ankle) improved percent agreement for activity type
354 classification over any single accelerometer but did not improve EE prediction over a
355 thigh-worn accelerometer in an adult sample. Additionally, studies examining the IDEEA

356 monitor, a five-accelerometer system (left and right upper leg, left and right foot,
357 sternum), generally show better EE prediction than some, but not all, single-
358 accelerometer prediction models in adults (Dannecker, Sazonova, Melanson, Sazonov, &
359 Browning, 2013; Lof, Henriksson, & Forsum, 2013; Ryan & Gormley, 2013). In contrast,
360 a previous study from our group (Mackintosh et al., 2016) demonstrated no additional EE
361 prediction accuracy in youth when combining up to eight combinations of accelerometer
362 locations relative to either hip or wrist alone in youth. Given that multi-accelerometer
363 systems may provide only small, if any, additional EE prediction accuracy, their utility
364 may be limited for EE prediction given the additional burden to researchers as well as
365 participants.

366 Our study has several notable strengths. Specifically, the direct comparison of two
367 popular accelerometer placement sites and their combination as well as count-based and
368 raw data offers important considerations for how to use accelerometers for EE prediction.
369 Additionally, the use of an independent sample for cross-validation is a strength of the
370 present study. Previous studies that aimed to develop machine learning models for EE
371 prediction have often used leave-one-out, *k*-fold, or other similar holdout
372 development/testing methods, which allow for training and testing of models to be
373 conducted efficiently within small samples. Because training and testing is being
374 conducted using data from the same study, there is inherent similarity in the types of
375 activities performed, setting, available equipment, and participant recruitment (Shao,
376 1993). Unsurprisingly, activity type or intensity prediction models developed in youth
377 and adult samples have yielded lower accuracy when evaluated in an independent cross-
378 validation (Gyllensten & Bonomi, 2011; Hibbing, Ellingson, et al., 2018; Kerr et al.,

379 2016; Sasaki et al., 2016). While our cross-validation sample had overlap in one of the
380 four exergames used compared to the calibration sample, there were still differences in
381 several of the activities (e.g., graded exercise test in calibration only), setting,
382 recruitment, and equipment that make the cross-validation sample independent from the
383 sample used to calibrate the models. However, future studies should aim to cross-validate
384 these models in independent samples participating in a larger variety of activities (e.g.,
385 ball games, tag).

386 Several study limitations must also be acknowledged. The present study was a
387 secondary analysis of data from a protocol conducted in a laboratory setting using only
388 exergames and, in calibration, a graded exercise test, resulting in limited activity types as
389 well as a high proportion of time spent in MVPA and low time spent sedentary.
390 Therefore, the applicability of the developed ANNs to other activity types, less active
391 portions of a youth's day, or for free-living EE prediction is unknown. Given the
392 overestimation of EE during rest/transitions by most ANNs and the regression models
393 during this study as well as the intercepts for the regression models falling close to 2.5
394 METs, these models are not suited to detecting time spent in sedentary behaviour and are,
395 therefore, only potentially useable for predicting MVPA . As such, these models will
396 have poorer accuracy if used to predict EE across a full day. Second, our comparison of
397 count-based and raw data, as well as different accelerometer placements, used only one
398 type of machine learning model, one set of features, and one epoch length. Performance
399 of different types of models, different feature sets, and across different epoch lengths may
400 yield informative results and should be explored in future research. Due to the more
401 varied movements at the wrist than the hip, it may be that more complex features of the

402 raw accelerometer data and/or other sensor inputs such as gyroscope, barometric
403 pressure, or heart rate may aid in EE prediction accuracy from wrist-worn accelerometers
404 (Wang et al., 2012). Finally, our raw accelerometer data were not autocalibrated, as is
405 recommended with raw data collection (van Hees et al., 2014), due to too short of a data
406 collection session in our laboratory setting. Autocalibration in a similar dataset from our
407 research team revealed calibration errors of ~2.2% which, although minor, could
408 potentially impact the machine learning models developed from raw data, so this should
409 be evaluated in future work.

410 **Conclusions**

411 In summary, our study found that a machine learning model developed from
412 count-based, hip accelerometer data had higher EE prediction accuracy in youth during
413 an exergame session than count-based models developed from wrist data or a
414 combination of hip and wrist data and higher accuracy than corresponding raw data
415 models and count-based regression equations. Although our results should be confirmed
416 using other types of machine learning models and feature sets and be expanded to include
417 activities other than exergames, our preliminary findings suggest that a hip-worn
418 accelerometer will provide better accuracy than wrist-worn accelerometers for EE
419 assessment in youth for assessing MVPA during exergames. On a separate note, we
420 recommend that transparent methods for filtering and processing raw accelerometer data
421 be developed to improve accuracy and comparability of accelerometer-based EE
422 prediction.

423 **Acknowledgements**

424 The authors would like to acknowledge David Barnett and Frazer Ashman as well
425 as funding from Bridging the Gaps for the Swansea portion of the study data and Wei
426 Peng, Jourdin Barkman, Allie Diltz, and funding from Michigan State University
427 Department of Media and Information for the Michigan State University portion of the
428 study data. Additionally, the authors would like to thank Jan Brønd and Daniel Arvidsson
429 for resampling the data from our calibration dataset so that it would be comparable to the
430 cross-validation dataset.

431 **References**

- 432 ActiGraph (2016). What are counts? [Online]. Available:
433 <https://actigraph.desk.com/customer/en/portal/articles/2515580-what-are-counts->
434 [Accessed July 2 2017].
- 435 Bai, J., DI, C., Xiao, L., Evenson, K. R., LaCroix, A. Z., Crainiceanu, C. M. & Buchner,
436 D. M. (2016). An activity index for raw accelerometry data and its comparison
437 with other activity metrics. *PLoS One*, 11, e0160644.
- 438 Bailey, R. C., Olson, J., Pepper, S. L., Porszasz, J., Barstow, T. J., & Cooper, D. M.
439 (1995). The level and tempo of children's physical activities: an observational
440 study. *Medicine and Science in Sports and Exercise*, 27(7), 1033-1041.
- 441 Bakrania, K., Yates, T., Rowlands, A. V., Esliger, D. W., Bunnell, S., Sanders, J.,
442 Davies, M., Khunti, K. & Edwardson, C. L. (2016). Intensity thresholds on raw
443 acceleration data: Euclidean norm minus one (ENMO) and mean amplitude
444 deviation (MAD) approaches. *PLoS One*, 11, e0164045.
- 445 Barkman, J., Pfeiffer, K.A., Diltz, A. & Peng, W. (2016). Examining energy expenditure
446 in youth Using Xbox Kinect: Differences by player mode. *Journal of Physical*
447 *Activity and Health*, 13, S41-3.
- 448 Bland, J. M., & Altman, D. G. (1986). Statistical methods for assessing agreement
449 between two methods of clinical measurement. *The Lancet*, 1(8476), 307-310.
- 450 Bouten, C. V., Sauren, A. A., Verduin, M., & Janssen, J. D. (1997). Effects of placement
451 and orientation of body-fixed accelerometers on the assessment of energy
452 expenditure during walking. *Medical & Biological Engineering & Computing*,
453 35(1), 50-56.
- 454 Brønd, J. C., Andersen, L. B., & Arvidsson, D. (2017). Generating ActiGraph counts
455 from raw acceleration recorded by an alternative monitor. *Medicine and Science*
456 *in Sports and Exercise*, 49(11), 2351-2360.
- 457 Brønd, J. C., & Arvidsson, D. (2016). Sampling frequency affects the processing of
458 ActiGraph raw acceleration data to activity counts. *Journal of Applied*
459 *Physiology*, 120(3), 362-369.

- 460 Clevenger, K. A. & Howe, C. A. (2015). Energy cost and enjoyment of active
461 videogames in children and teens: Xbox 360 Kinect. *Games for Health Journal*,
462 4, 318-24.
- 463 Crouter, S. E., Flynn, J. I., & Bassett, D. R., Jr. (2015). Estimating physical activity in
464 youth using a wrist accelerometer. *Medicine and Science in Sports and Exercise*,
465 47(5), 944-951.
- 466 Crouter, S. E., Horton, M., & Bassett, D. R., Jr. (2012). Use of a two-regression model
467 for estimating energy expenditure in children. *Medicine and Science in Sports and*
468 *Exercise*, 44(6), 1177-1185.
- 469 Crouter, S. E., Oody, J. F., & Bassett, D. R., Jr. (2018). Estimating physical activity in
470 youth using an ankle accelerometer. *Journal of Sports Sciences*, 36(19), 2265-
471 2271.
- 472 Dannecker, K. L., Sazonova, N. A., Melanson, E. L., Sazonov, E. S. & Browning, R. C.
473 (2013). A comparison of energy expenditure estimation of several physical
474 activity monitors. *Medicine and Science in Sports and Exercise*, 45, 2105-12.
- 475 Dong, B., Montoye, A., Moore, R., Pfeiffer, K. & Biswas, S. (2013). Energy-aware
476 activity classification using wearable sensor networks. *Proceedings of the SPIE*
477 *International Society for Optical Engineering*, 2013, 87230Y-87230Y.
- 478 Dong, B., Biswas, S., Montoye, A., & Pfeiffer, K. (2013). Comparing metabolic energy
479 expenditure estimation using wearable multi-sensor network and single
480 accelerometer. *Conference proceedings of the IEEE Engineering in Medicine and*
481 *Biology Society*, 2013, 2866-2869.
- 482 Ellis, K., Kerr, J., Godbole, S., Lanckriet, G., Wing, D., & Marshall, S. (2014). A random
483 forest classifier for the prediction of energy expenditure and type of physical
484 activity from wrist and hip accelerometers. *Physiological Measurement*, 35(11),
485 2191-2203.
- 486 Esteban-Cornejo, I., Tejero-Gonzalez, C. M., Sallis, J. F. & Veiga, O. L. (2015). Physical
487 activity and cognition in adolescents: A systematic review. *Journal of Science and*
488 *Medicine in Sport*, 18, 534-9.
- 489 FAO/WHO/UNU Expert Consultation (2001). Energy requirements of adults. Retrieved
490 from <http://www.fao.org/docrep/007/y5686e/y5686e07.htm>
- 491 Freedson, P., Pober, D. & Janz, K. F. (2005). Calibration of accelerometer output for
492 children. *Medicine and Science in Sports and Exercise*, 37, S523-30.
- 493 Glickman, M. E., Rao, S. R., & Schultz, M. R. (2014). False discovery rate control is a
494 recommended alternative to Bonferroni-type adjustments in health studies.
495 *Journal of Clinical Epidemiology*, 67(8), 850-857.
- 496 Graves, L. E., Ridgers, N. D. & Stratton, G. (2008). The contribution of upper limb and
497 total body movement to adolescents' energy expenditure whilst playing Nintendo
498 Wii. *European Journal of Applied Physiology*, 104, 617-23.
- 499 Gyllensten, I. C. & Bonomi, A. G. (2011). Identifying types of physical activity with a
500 single accelerometer: Evaluating laboratory-trained algorithms in daily life. *IEEE*
501 *Transactions on Biomedical Engineering*, 58, 2656-63.
- 502 Hibbing, P. R., Ellingson, L. D., Dixon, P. M., & Welk, G. J. (2018). Adapted Sojourn
503 Models to Estimate Activity Intensity in Youth: A Suite of Tools. *Medicine and*
504 *Science in Sports and Exercise*, 50(4), 846-854.

- 505 Hibbing, P. R., LaMunion, S. R., Kaplan, A. S., & Crouter, S. E. (2018). Estimating
506 energy expenditure with ActiGraph GT9X inertial measurement unit. *Medicine
507 and Science in Sports and Exercise*, 50(5), 1093-1102. Hwang, J., Fernandez, A.
508 M., & Lu, A. S. (2018). Application and validation of activity monitors' epoch
509 lengths and placement sites for physical activity assessment in exergaming.
510 *Journal of Clinical Medicine*, 7(9).
- 511 John, D. & Freedson, P. (2012). ActiGraph and Actical physical activity monitors: A
512 peek under the hood. *Medicine and Science in Sports and Exercise*, 44, S86-9.
- 513 Kerr, J., Patterson, R. E., Ellis, K., Godbole, S., Johnson, E., Lanckriet, G., &
514 Staudenmayer, J. (2016). Objective assessment of physical activity: Classifiers for
515 public health. *Medicine and Science in Sports and Exercise*, 48(5), 951-957.
- 516 Lof, M., Henriksson, H. & Forsum, E. (2013). Evaluations of Actiheart, IDEEA and RT3
517 monitors for estimating activity energy expenditure in free-living women. *Journal
518 of Nutritional Sciences*, 2, e31.
- 519 Lyden, K., Keadle, S. K., Staudenmayer, J. & Freedson, P. S. (2014). A method to
520 estimate free-living active and sedentary behavior from an accelerometer.
521 *Medicine and Science in Sports and Exercise*, 46, 386-97.
- 522 Lyons, R. (2013). *Understanding Digital Signal Processing* (3rd ed.). Upper Saddle
523 River, NJ: Prentice Hall.
- 524 Mackintosh, K. A., Montoye, A. H. K., Pfeiffer, K. A. & McNarry, M. (2016).
525 Investigating optimal accelerometer placement for energy expenditure prediction
526 in children using a machine learning approach. *Physiological Measurement*, 37,
527 1728-1740.
- 528 McManus, C. (1991). The inheritance of left-handedness. In: SYMPOSIUM, C. F. (ed.)
529 *Biological Asymmetry and Handedness*. West Sussex, England: John Wiley &
530 Sons Ltd.
- 531 McMurray, R. G., Butte, N. F., Crouter, S. E., Trost, S. G., Pfeiffer, K. A., Bassett, D. R.,
532 Puyau, M. R., Berrigan, D., Watson, K. B., & Fulton, J. E. (2015). Exploring
533 metrics to express energy expenditure of physical activity in youth. *PLoS One*, 10,
534 e0130869.
- 535 Montoye, A. H. K., Begum, M., Henning, Z. & Pfeiffer, K. A. (2017). Comparison of
536 linear and non-linear models for predicting energy expenditure from raw
537 accelerometer data. *Physiological Measurement*, 38, 343-357.
- 538 Montoye, A. H. K., Moore, R. W., Bowles, H. R., Korycinski, R. & Pfeiffer, K. A.
539 (2016). Reporting accelerometer methods in physical activity intervention studies:
540 a systematic review and recommendations for authors. *British Journal of Sports
541 Medicine*.
- 542 Montoye, A. H. K., Mudd, L. M., Biswas, S. & Pfeiffer, K. A. (2015). Energy
543 expenditure prediction using raw accelerometer data in simulated free living.
544 *Medicine and Science in Sports and Exercise*, 47, 1735-46.
- 545 Montoye, A. H. K., Pivarnik, J. M., Mudd, L. M., Biswas, S. & Pfeiffer, K. A. (2016).
546 Wrist-independent energy expenditure prediction models from raw accelerometer
547 data. *Physiological Measurement*, 37, 1770-1784.
- 548 O'Driscoll, R., Turicchi, J., Beaulieu, K., Scott, S., Matu, J., Deighton, K., . . . Stubbs, J.
549 (2018). How well do activity monitors estimate energy expenditure? A systematic

- 550 review and meta-analysis of the validity of current technologies. *British Journal*
551 *of Sports Medicine*.
- 552 Preece, S. J., Goulermas, J. Y., Kenney, L. P., Howard, D., Meijer, K., & Crompton, R.
553 (2009). Activity identification using body-mounted sensors: A review of
554 classification techniques. *Physiological Measurement*, 30(4), R1-33.
555 doi:10.1088/0967-3334/30/4/R01
- 556 Ripley, B., & Venables, W. (2016). Package 'nnet': Feed-Forward Neural Networks and
557 Multinomial Log-Linear Models. Retrieved from [https://cran.r-](https://cran.r-project.org/web/packages/nnet/nnet.pdf)
558 [project.org/web/packages/nnet/nnet.pdf](https://cran.r-project.org/web/packages/nnet/nnet.pdf)
- 559 Rosdahl, H., Gullstrand, L., Salier-Eriksson, J., Johansson, P., & Schantz, P. (2010).
560 Evaluation of the Oxycon Mobile metabolic system against the Douglas bag
561 method. *European Journal of Applied Physiology*, 109(2), 159-171.
- 562 Rosenberg, M., Lay, B., Lee, M., Derbyshire, A., Kur, J., Ferguson, R., Maitland, C.,
563 Mills, A., Davies, C., Pratt, I. S. & Braham, R. (2013). New-generation active
564 videogaming maintains energy expenditure in children across repeated bouts.
565 *Games for Health Journal*, 2, 274-9.
- 566 Ryan, J. & Gormley, J. (2013). An evaluation of energy expenditure estimation by three
567 activity monitors. *European Journal of Sports Science*, 13, 681-8.
- 568 Sasaki, J. E., John, D., & Freedson, P. S. (2011). Validation and comparison of
569 ActiGraph activity monitors. *Journal of Science and Medicine in Sport*, 14(5),
570 411-416.
- 571 Sasaki, J. E., Hickey, A. M., Staudenmayer, J. W., John, D., Kent, J. A. & Freedson, P. S.
572 (2016). Performance of activity classification algorithms in free-living older
573 adults. *Medicine and Science in Sports and Exercise*, 48, 941-50.
- 574 Schofield, W. N. (1985). Predicting basal metabolic rate, new standards and review of
575 previous work. *Human Nutrition and Clinical Nutrition*, 39 Suppl 1, 5-41.
- 576 Shao, J. (1993). Linear model selection by cross-validation. *Journal of the American*
577 *Statistical Association*, 88, 486-494.
- 578 Staudenmayer, J., He, S., Hickey, A., Sasaki, J. & Freedson, P. (2015). Methods to
579 estimate aspects of physical activity and sedentary behavior from high-frequency
580 wrist accelerometer measurements. *Journal of Applied Physiology*, 119, 396-403.
- 581 Staudenmayer, J., Pober, D., Crouter, S., Bassett, D. & Freedson, P. (2009). An artificial
582 neural network to estimate physical activity energy expenditure and identify
583 physical activity type from an accelerometer. *Journal of Applied Physiology*, 107,
584 1300-7.
- 585 Swartz, A. M., Strath, S. J., Bassett, D. R., Jr., O'Brien, W. L., King, G. A. & Ainsworth,
586 B. E. (2000). Estimation of energy expenditure using CSA accelerometers at hip
587 and wrist sites. *Medicine and Science in Sports and Exercise*, 32, S450-6.
- 588 Troiano, R. P., Berrigan, D., Dodd, K. W., Masse, L. C., Tilert, T. & McDowell, M.
589 (2008). Physical activity in the United States measured by accelerometer.
590 *Medicine and Science in Sports and Exercise*, 40, 181-8.
- 591 Troiano, R. P., McClain, J. J., Brychta, R. J. & Chen, K. Y. (2014). Evolution of
592 accelerometer methods for physical activity research. *British Journal of Sports*
593 *Medicine*, 48, 1019-1023.

- 594 Trost, S. G., Wong, W. K., Pfeiffer, K. A. & Zheng, Y. (2012). Artificial neural networks
595 to predict activity type and energy expenditure in youth. *Medicine and Science in*
596 *Sports and Exercise*, 44, 1801-9.
- 597 US Department of Health and Human Services. (2018). Physical Activity Guidelines
598 Advisory Committee: 2018 Physical Activity Guidelines for Americans, 2nd
599 edition. Available: www.health.gov/paguidelines.
- 600 van Hees, V. T., Fang, Z., Langford, J., Assah, F., Mohammad, A., da Silva, I. C., . . .
601 Brage, S. (2014). Autocalibration of accelerometer data for free-living physical
602 activity assessment using local gravity and temperature: an evaluation on four
603 continents. *Journal of Applied Physiology*, 117(7), 738-744.
- 604 van Hees, V. T., Gorzelniak, L., Dean Leon, E. C., Eder, M., Pias, M., Taherian, S., . . .
605 Brage, S. (2013). Separating movement and gravity components in an
606 acceleration signal and implications for the assessment of human daily physical
607 activity. *PLoS One*, 8(4), e61691.
- 608 van Hees, V. T., Sabia, S., Anderson, K. N., Denton, S. J., Oliver, J., Catt, M., . . . Singh-
609 Manoux, A. (2015). A novel, open access method to assess sleep duration using a
610 wrist-worn accelerometer. *PLoS One*, 10(11), e0142533.
- 611 Wang, J., Redmond, S. J., Voleno, M., Narayanan, M. R., Wang, N., Cerutti, S., &
612 Lovell, N. H. (2012). Energy expenditure estimation during normal ambulation
613 using triaxial accelerometry and barometric pressure. *Physiological Measurement*,
614 33(11), 1811-1830.
- 615

616 **Tables**617 Table 1. Participant characteristics in calibration and cross-validation samples.
618

	Calibration			Cross-validation		
	Boys (n=15)	Girls (n=12)	Total (n=27)	Boys (n=21)	Girls (n=13)	Total (n=34)
Age (years)	10.8 (1.2)	10.8 (1.4)	10.8 (1.0)	11.6 (1.4)	11.5 (1.1)	11.6 (1.2)
Stature (m)	1.46 (0.13)	1.45 (0.10)	1.45 (0.11)	1.57 (0.14)	1.54 (0.11)	1.56 (0.12)
Mass (kg)	38.7 (8.5)	37.2 (9.0)	38.7 (8.8)	49.3 (14.2)	48.1 (15.1)	48.8 (14.6)
Predicted basal VO₂ (ml·kg⁻¹·min⁻¹)	4.9 (0.4)	4.6 (0.7)	4.8 (0.5)	4.4 (0.5)	4.0 (0.7)	4.3 (0.6)

619 Data are shown as mean (standard deviation).
620
621

Table 2. Criterion-measured and accelerometer-predicted energy expenditure in cross-validation.

	Overall (n=34)	Rest/ transition (n=34)	Kinect			Kinect Sports Boxing (n=18)
			Reflex Ridge (n=17)	Just Dance 3 (n=17)	Wipeout (n=16)	
Criterion	3.9 (1.5)	2.7 (1.1)	4.4 (1.5)	4.4 (1.5)	4.0 (1.5)	3.2 (1.0)
Hip count ANN	4.0 (1.2)	3.3 (0.7)*	4.6 (1.4)	4.4 (1.2)	3.7 (0.8)*	3.7 (0.8)*
Hip raw ANN	3.0 (1.6)*	2.7 (1.4)	3.1 (1.9)*	3.3 (1.5)*	2.9 (1.7)*	2.9 (1.4)*
Wrist count ANN	5.3 (1.6)*	4.1 (1.2)*	4.9 (1.2)*	5.3 (1.3)*	5.2 (1.4)*	6.5 (1.7)*
Wrist raw ANN	5.9 (3.7)*	3.8 (2.8)*	4.8 (3.1)*	6.2 (3.4)*	7.6 (4.5)*	5.6 (2.9)*
Combination count ANN	4.7 (1.7)	3.5 (1.0)*	4.9 (1.6)*	5.1 (1.7)*	4.4 (1.6)*	4.9 (1.6)*
Combination raw ANN	3.3 (2.0)	2.8 (1.7)	3.6 (2.5)*	3.3 (1.6)*	3.9 (2.3)	2.7 (1.1)*
Regression hip VM	4.3 (1.7)	3.2 (1.1)*	5.1 (2.2)*	4.8 (1.8)*	4.0 (1.3)	3.7 (1.0)*
Regression wrist VM	6.1 (2.3)*	4.3 (1.7)*	5.4 (1.7)*	5.8 (1.9)*	6.7 (2.9)*	7.3 (2.1)*
Regression combination VM	5.5 (2.0)*	3.8 (1.4)*	5.5 (1.8)*	5.4 (1.7)*	5.4 (2.0)*	5.4 (1.5)*

Data are shown in metabolic equivalents (METs), as mean (standard deviation).

*Indicates significant difference from the criterion.

Combination: Combination of hip and wrist data.

ANN: Artificial neural network machine learning model.

VM: regression equation developed using vector magnitude of count-based data.

Figure titles and captions

Figure 1. Root mean squared error for energy expenditure prediction.

¹Indicates significant difference from all other models.

²Indicates significant difference from all models except ANN wrist count and ANN combination raw models.

³Indicates significant difference from all models except ANN hip raw and ANN combination raw models.

⁴Indicates significant difference from all models except regression wrist model.

⁵Indicates significant difference from all models except ANN combination raw model.

⁶Indicates significant difference from ANN hip count, ANN wrist raw, regression wrist, and regression combination models.

⁷Indicates significant difference from all models except ANN combination raw model.

⁸Indicates significant difference from all models except ANN wrist raw model.

ANN: artificial neural network.

Combo: Combination of hip and wrist data.

METs: Metabolic equivalents.

Figure 2. Bland-Altman plots showing agreement between predicted and measured energy expenditure when cross-validating artificial neural networks and regression models.

- a. ANN developed from count-based, hip accelerometer data.
- b. ANN developed from raw, hip accelerometer data.
- c. ANN developed from count-based, wrist accelerometer data.
- d. ANN developed from raw, wrist accelerometer data.
- e. ANN developed from count-based, combination accelerometer data.
- f. ANN developed from raw, combination accelerometer data.
- g. Regression model developed from count-based, hip accelerometer data.
- h. Regression model developed from count-based, wrist accelerometer data.
- i. Regression model from count-based, combination accelerometer data.

ANN: artificial neural network.

Combo: Combination of hip and wrist data.

METs: Metabolic equivalents.

Points greater than 0 on the y-axis represent underestimation by the predictive model, and vice versa for points less than 0.