



**UNIVERSITY OF LEEDS**

This is a repository copy of *Improving energy expenditure estimates from wearable devices: A machine learning approach*.

White Rose Research Online URL for this paper:  
<https://eprints.whiterose.ac.uk/159856/>

Version: Accepted Version

---

**Article:**

O'Driscoll, R, Turicchi, J [orcid.org/0000-0003-1174-813X](https://orcid.org/0000-0003-1174-813X), Hopkins, M [orcid.org/0000-0002-7655-0215](https://orcid.org/0000-0002-7655-0215) et al. (3 more authors) (2020) Improving energy expenditure estimates from wearable devices: A machine learning approach. *Journal of Sports Sciences*, 38 (13). pp. 1496-1505. ISSN 0264-0414

<https://doi.org/10.1080/02640414.2020.1746088>

---

© 2020 Informa UK Limited, trading as Taylor & Francis Group. This is an author produced version of a journal article published in *Journal of Sports Sciences*. Uploaded in accordance with the publisher's self-archiving policy.

**Reuse**

Items deposited in White Rose Research Online are protected by copyright, with all rights reserved unless indicated otherwise. They may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



[eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk)  
<https://eprints.whiterose.ac.uk/>

1 Improving Energy Expenditure Estimates From Wearable Devices: A Machine Learning  
2 Approach

3

4 O'Driscoll Ruairi<sup>1</sup>, Turicchi Jake<sup>1</sup>, Hopkins Mark<sup>2</sup>, Horgan Graham. W.<sup>3</sup>, Finlayson Graham<sup>1</sup>,  
5 Stubbs James. R.<sup>1</sup>.

6 <sup>1</sup>*Appetite Control and Energy Balance Group, School of Psychology, University of Leeds,*  
7 *Leeds, U.K*

8 <sup>2</sup>*School of Food Science and Nutrition, Faculty of Mathematics and Physical Sciences,*  
9 *University of Leeds, Leeds, U.K.*

10 <sup>3</sup>*Biomathematics & Statistics Scotland, Aberdeen, U.K.*

11 *Email addresses:*

12 O'Driscoll Ruairi ([psrod@leeds.ac.uk](mailto:psrod@leeds.ac.uk)), Turicchi Jake ([psjt@leeds.ac.uk](mailto:psjt@leeds.ac.uk)), Hopkins Mark  
13 ([M.Hopkins@leeds.ac.uk](mailto:M.Hopkins@leeds.ac.uk)), Horgan Graham. W. ([graham.horgan@bioss.ac.uk](mailto:graham.horgan@bioss.ac.uk)), Finlayson  
14 Graham ([g.s.finlayson@leeds.ac.uk](mailto:g.s.finlayson@leeds.ac.uk)), Stubbs James. R. ([r.j.stubbs@leeds.ac.uk](mailto:r.j.stubbs@leeds.ac.uk)).

15 *Keywords:*

16 **Machine learning,**

17 **Heart rate,**

18 **Energy Expenditure,**

19 **Accelerometer**

**20 Abstract**

21 A means of quantifying continuous, free-living energy expenditure (EE) would  
22 advance the study of bioenergetics. The aim of this study was to apply a non-linear, machine  
23 learning algorithm (random forest) to predict minute level EE for a range of activities using  
24 acceleration, physiological signals (e.g. heart rate, body temperature, galvanic skin response),  
25 and participant characteristics (e.g. sex, age, height, weight, body composition) collected  
26 from wearable devices (Fitbit charge 2, Polar H7, SenseWear Armband Mini and Actigraph  
27 GT3-x) as potential inputs. By utilising a leave-one-out cross-validation approach in 59  
28 subjects, we investigated the predictive accuracy in sedentary, ambulatory, household, and  
29 cycling activities compared to indirect calorimetry (Vyntus CPX). Over all activities,  
30 correlations of at least  $r=0.85$  were achieved by the models. Root mean squared error ranged  
31 from 1-1.37 METs and all overall models were statistically equivalent to the criterion  
32 measure. Significantly lower error was observed for Actigraph and Sensewear models, when  
33 compared to the manufacturer provided estimates of the Sensewear Armband ( $p<0.05$ ). A  
34 high degree of accuracy in EE estimation was achieved by applying non-linear models to  
35 wearable devices which may offer a means to capture the energy cost of free-living activities.

## 36 **Background**

37           The measurement of energy expenditure (EE) is critical to understand human energy  
38 requirements in health and disease, and how components of energy balance contribute to over  
39 and under nutrition. Quantifying energy balance in free-living individuals requires the precise  
40 and accurate estimation of at least two of the three components of energy balance; energy  
41 intake (EI), EE and changes in energy stored (ES). Currently, quantifying true patterns of EI  
42 and EE in the free-living environment is constrained by methodological and practical  
43 limitations. Objective measures of EI can be derived from EE and changes in ES (e.g. the  
44 intake-balance method [1]), but the use of doubly labelled water (DLW) to estimate EE is  
45 costly and fails to capture daily variation in EE, which limits its widespread adoption [2].

46           Activity monitors have long been recognised as a potential means to estimate EE [3],  
47 but the accuracy of current devices in estimating the energy cost of a wide range of activities  
48 and intensities is limited [4]. Accelerometry is routinely used to quantify bodily movement  
49 and to predict EE using linear models [5,6]. However, the relationship between EE and  
50 acceleration is variable between activities [7], and accelerometry alone has limited  
51 sensitivity to capture the additional energy demands of activities that do not alter the velocity  
52 of movement (e.g. load carrying or incline walking) [8]. While estimates of EE from devices  
53 with HR sensing technology are likely to have reduced error relative to devices based on  
54 accelerometry alone [4], the relationship between accelerometry, HR and EE exhibits  
55 linearity only within specific activity types. Combined linear models that estimate EE will  
56 therefore not generalise across the range of human activities [9]. In some cases, this may  
57 explain the demonstrably poor manufacturer provided EE estimates obtained from many  
58 current and past activity monitors [10].

59           Complex, non-linear machine learning algorithms applied to research-grade  
60 accelerometers have shown remarkable accuracy in estimating EE using tree-based methods

61 [11] and artificial neural networks [9,12] and they may offer a means to overcome the limited  
62 accuracy of current devices. However, whether machine learning can be used to improve the  
63 estimation of EE using the sensor data obtained from commercial devices has yet to be  
64 examined.

65 Machine learning methods demonstrate clear potential to estimate EE using data  
66 obtained from wearable sensors. This study aims to explore the potential for non-linear,  
67 machine learning regression models utilising subject characteristics, movement and  
68 physiological variables to estimate EE in a range of activities.

69

## 70 **Methods**

### 71 **Participants**

72 A sample of 59 participants were included (Female=41, Age =  $44.4 \pm 14.1$  years,  
73 Weight =  $75.7 \pm 13.6$  kg, BMI =  $26.9 \pm 4.7$  kg/m<sup>2</sup>, FM=  $24.8 \pm 10.73$  kg, FFM =  $49.8 \pm 8.9$ ,  
74 FM (%) =  $32.5 \pm 10.3\%$ , FFM (%) =  $67.5 \pm 10.3\%$ , RMR =  $1581.8 \pm 280.4$  kcal/d).

75 Participants were primarily from the Leeds centre of the NoHoW trial (n = 44), ISRCTN  
76 registry (ISRCTN88405328), an additional 15 participants were recruited from the University  
77 of Leeds and surrounding areas. Participants were excluded from the study for the following  
78 reasons: pregnancy, medications altering metabolic rate, cardiovascular, metabolic, renal  
79 disorders, illness or injury that provide an increased risk of medical events during PA [13].  
80 Ethical approval was granted by The University of Leeds, School of Psychology ethics  
81 committee (PSC-407, 18/08/2018).

### 82 **Physical measurements**

83 Measurements were conducted at an exercise laboratory at The Human Appetite  
84 Research Unit, University of Leeds. Participants arrived at the university between 06:00 am  
85 and 09:30 am, having refrained from the intake of food, caffeine and exercise for 12 hours

86 prior to the measurements. Systolic and diastolic blood pressure (BP) and resting heart rate  
87 (HR) (Microlife BP A2 Basic, Gentle Technology, Microlife, Clearwater, FL, USA, Inc.)  
88 were measured at rest and in the sitting position. Height ( $\pm 0.1$ cm) was measured barefoot,  
89 using a Seca 704s instrument (SECA, Germany). Fat mass (FM) and fat-free mass (FFM)  
90 were estimated using air displacement plethysmography (BodPod, Life Measurement, Inc.;  
91 USA) and the Siri equation [14]. Body weight ( $\pm 0.1$ kg) was also obtained from the BodPod  
92 scales in light clothing.

93 **Table 1 – insert here**

94 **Physical activity devices**

95 Participants wore a number of physical activity devices during the study and all  
96 devices were initialised in accordance with manufacturer's instructions. The Polar m400 HR  
97 Monitor Watch and a Polar H7 chest strap (Polar Electro, Kempele, Finland) were used to  
98 measure HR. The Polar H7 facilitates approximately 200 hours of continuous transmission.  
99 In this study data were extracted at the second level and averaged to the minute-level.  
100 Participants also wore a Fitbit Charge 2 (FC2) (Fitbit Inc, San Francisco, CA, USA), a wrist-  
101 worn activity monitor, which incorporates a tri-axial accelerometer. The FC2 also estimates  
102 HR through a patented technology called 'PurePulse', which uses light-emitting diodes to  
103 measure changes to blood volume [15]. An Actigraph GT3-x accelerometer (AG) was placed  
104 on the non-dominant wrist which measured acceleration along vertical, horizontal and  
105 perpendicular axes at a sample rate of 30Hz. Participants also wore the SenseWear Armband  
106 Mini (SWA) (BodyMedia Inc., Pittsburgh, PA) on the non-dominant upper arm. The SWA  
107 collected tri-axial accelerometer data and as well as data from heat-related sensors measuring  
108 heat flux, skin temperature, near body ambient temperature and galvanic skin response.

109 **Energy expenditure measurement**

110 Resting metabolic rate (RMR) was measured using an indirect calorimetry system  
111 fitted with a ventilated hood (GEM, Nutren Technology Ltd; UK). Participants lay in the  
112 supine position for 30 minutes, whilst  $\text{VO}_2$  and  $\text{VCO}_2$  were continually measured. An RMR  
113 estimate was derived from 5 minutes of steady state data, as described previously [16].  
114 Briefly, after discarding the first 5 minutes, the  $\text{VO}_2$  and  $\text{VCO}_2$  measurements in the 5-minute  
115 period with the lowest coefficient of variation during the overall measurement period are  
116 used to estimate RMR. In the absence of an RMR measurement ( $n=2$ ), a body mass index  
117 (BMI) specific RMR algorithm was used [17]. A stationary metabolic cart (Vyntus CPX,  
118 Jaeger-CareFusion, UK) was used as the criterion measure of EE in the physical activity  
119 protocol. Breath by breath data were aggregated to minute level to estimate EE ( $\text{kcal}\cdot\text{min}^{-1}$ ).  
120 The Vyntus CPX is highly valid and reliable [18,19] and served as a criterion comparison for  
121 the developed models. The unit was calibrated prior to each lab visit in accordance with  
122 manufacturer's instructions. Data were aggregated to the minute level and EE ( $\text{kcal}\cdot\text{min}^{-1}$ )  
123 values were calculated from  $\text{VO}_2$  and  $\text{VCO}_2$  data assuming a minimal contribution of protein  
124 oxidation [20]. We expressed minute level EE as a multiple of each participant's RMR, to  
125 derive metabolic equivalents (METs), which was the outcome variable.

### 126 **Physical activity protocol**

127 Participants undertook a structured protocol consisting of 10 activities, which were  
128 performed at a consistent intensity of 5 minutes each. The activities were performed in a set  
129 order and included: sitting, standing, treadmill walking (4 km/h), incline treadmill walking (4  
130 km/h, 5% incline), jogging (6-8 km/h, 5% incline), incline jogging (6-8 km/h, 5% incline).  
131 Next, after a 3-minute resting period, participants transitioned to a cycle ergometer for low-  
132 intensity (30 watts), and moderate intensity cycling (60 watts). After another period of  
133 recovery, participants performed a folding task and lastly a sweeping task. The physical  
134 activity protocol was performed by all participants, however the jogging task ( $n=49$ ), the

135 jogging 5% incline (n=30) and the moderate cycling tasks (n=58) were not performed by all  
136 participants, due to variation in physical fitness.

### 137 **Data processing and model development**

138 All data sources were aggregated to the minute-level and were matched by time for  
139 each participant. The first minute of data was removed leaving minutes 2-5 for inclusion in  
140 model development [11,21]. The models developed in this study were trained using complete  
141 minute level data only, so if any data points were missing from the sensors or subject  
142 characteristics, that single minute was not included in the analysis. In the present study, we  
143 developed distinct predictive models for each device (FC2, AG and SWA) and the specific  
144 predictor variables used in each of these is described in table 1. The algorithm used in the  
145 present study was a random forest regressor [22]. Random forests are an ensemble method  
146 which aggregate the output of numerous decision trees to produce a continuous output. In the  
147 random forest algorithm, trees are trained on a random sample of the available predictor  
148 variables, which reduces the chances overfitting the training data [23]. For all the random  
149 forest models, the number of variables randomly sampled at each split was set 1/3 of the  
150 number of predictor variables in the model. For each of the developed models, 1000 trees  
151 were grown, and minimum size of terminal nodes was set to 5. All model development and  
152 training was conducted with the “randomForest” package [24] in R. Model parameters were  
153 established in preliminary tuning experiments and were standardised to allow comparability  
154 between each of the models.

### 155 **Statistical analysis**

156 Two validation approaches were used in this study. A ‘holdout’ approach was used in  
157 which all available data are split into training and testing, at a ratio of 80:20. Secondly,  
158 Leave-One-Out Cross-Validation (LOOCV), in which models are trained on all participants’



159 data with the exception of one participant, which serves as the testing dataset. This process is  
160 repeated until all participants data has been used to test the algorithm.

161

162 A range of statistical tests were employed to investigate the accuracy of model  
163 estimates, in line with previous validation research [25]: Pearson’s correlation coefficients,  
164 root mean squared error (RMSE) and mean absolute percentage error (MAPE), calculated  
165 with the R package ‘metrics’ [26]. Equivalence tests were used to determine whether the  
166 models were statistically equivalent to the criterion measured METs, to be considered  
167 equivalent, the 90% confidence interval of the estimate must fall within  $\pm 10\%$  of the  
168 criterion mean [27]. Repeated measures analysis of variance (ANOVA) tests were employed  
169 to test for differences in MAPE calculated for each of the models, and the SWA for each  
170 subject’s activity modality. We investigated differences between specific models with  
171 pairwise t-tests conducted with a Holm–Bonferroni false error rate correction. All data are  
172 reported as means and standard deviations (SD) unless otherwise stated. In order to estimate  
173 the precision of estimates in this study, standard errors for the overall RMSE of each model  
174 have been computed at the participant level and are presented in supplementary table 1.

175 All analyses and data processing were conducted R version 3.5.1 and RStudio  
176 Version 1.1.447 [28], using a p-value of  $< 0.05$  to determine statistical significance.

177

## 178 **Results**

### 179 **Predictive accuracy of models**

180 The performance of the holdout validation was typically superior to the LOOCV  
181 method, as measured by MAPE. The correlation of all models exceeded r values of 0.94 and  
182 the results of the models using the holdout approach are shown in supplementary table 2.

183 Models FBRF<sub>1</sub> and FBRF<sub>2</sub> demonstrated the greatest MAPE and highest RMSE using this  
184 approach.

185 **Table 2 – insert here**

186 The performance of the models without body composition data (AGRF<sub>2</sub>, SWRF<sub>2</sub>,  
187 FBRF<sub>1</sub>) using a LOOCV validation approach are shown in the form of scatterplots in figure  
188 1. The accuracy statistics from the LOOCV validation and the results of the equivalence tests  
189 are presented in table 2. Data loss occurred for 2 participant's FC2 data, one participant's AG  
190 data and one participant's polar HR data. All models were validated on at least 2000 minutes  
191 and 55 participants and individual level data is presented in supplementary table 3. The SWA  
192 was not statistically equivalent to the criterion measure, in contrast to the random forest  
193 models, which were all statistically equivalent to the criterion measure. The SWA also had  
194 the highest RMSE of 1.8 METs compared to the AGRF and SWRF models which ranged  
195 between 1-1.24 METs and FBRF models, which had RMSE values of 1.37 METs or less.

196 **Figure 1 – insert here**

197 **Table 2 – insert here**

198 The results of the ANOVA demonstrated a significant F statistic of 41.79 ( $p=7.26^{-49}$ )  
199 for between model differences, indicating that differences existed between the MAPE values  
200 for each model's METs estimates relative the criterion METs. Pairwise t-tests demonstrated  
201 that the MAPE for the SWA estimates were significantly higher than all random forest  
202 models ( $p<0.05$ ), except for FBRF models. Model AGRF<sub>1</sub> had MAPE values significantly  
203 higher than AGRF<sub>2</sub> ( $p=1.16^{-10}$ ) and AGRF<sub>3</sub> ( $p=2.03^{-08}$ ). SWRF<sub>1</sub> was significantly higher than  
204 SWRF<sub>2</sub> ( $p=7.19^{-11}$ ) and SWRF<sub>3</sub> ( $p=4.14^{-09}$ ).

205 The introduction of body composition did not result in significantly different MAPE  
206 models developed on AG sensor outputs, however, FBRF<sub>1</sub> had a significantly lower MAPE  
207 than the body composition model, FBRF<sub>2</sub> ( $p=0.007$ ) and SWRF<sub>2</sub> was significantly lower than

208 SWRF<sub>3</sub> (p=0.021), indicating a less accurate model performance with the addition of body  
209 composition. All SWRF and AGRF models had significantly lower MAPE values than  
210 FBRF<sub>1</sub> and FBRF<sub>2</sub> (p<0.01).

### 211 **Activity specific accuracy**

212 Activity specific accuracy statistics calculated using LOOCV are presented in table 3.  
213 The accuracy of the FBRF models was poorest in sedentary tasks, where MAPE values of  
214 47.87 and 52.30 were observed for FBRF<sub>1</sub> and FBRF<sub>2</sub>, respectively during the standing task.  
215 Both FBRF<sub>1</sub> and FBRF<sub>2</sub> were statistically equivalent during the jogging task. Models AGRF<sub>2</sub>  
216 and AGRF<sub>3</sub> were statistically equivalent in 5 of 10 tasks, namely: standing, incline walking,  
217 jogging, low intensity cycling and folding, and the AGRF<sub>1</sub> model was equivalent in all of the  
218 aforementioned tasks, except from incline walking. The MAPE values ranged from 13.01  
219 (AGRF<sub>3</sub>, walk incline) to 29.33 (AGRF<sub>1</sub>, sweeping). Models developed on SWA data were  
220 statistically equivalent for walking (SWRF<sub>2</sub> and SWRF<sub>3</sub> only), walking incline, jogging and  
221 low intensity cycling, and the SWRF<sub>2</sub> was equivalent in sitting. SWRF<sub>2</sub> and SWRF<sub>3</sub>  
222 demonstrated the poorest accuracy in the household tasks with overestimates in models and  
223 MAPE values ranging from 24.11 to 33.41.

### 224 **Table 3 – insert here**

### 225 **Model characteristics**

226 Using the feature set which included body composition data for each device (AGRF<sub>3</sub>,  
227 FBRF<sub>2</sub>, SWRF<sub>3</sub>), we computed the relative importance of each predictive variables [22]. The  
228 variable in the plots represents the percentage increase to the mean squared error following  
229 the permutation (random shuffling) of each variable. Permutation in this manner breaks the  
230 association between the predictive and outcome variable relative to the original model and  
231 therefore facilitates estimates of the importance of this variable to overall accuracy of the  
232 original model. Outlined in figure 2 for the FB, acceleration (i.e. steps) and HR normalised to

233 the sitting HR were the most important variables in the models, with age, kilograms of FFM,  
234 height and FM following after. In the SWRF<sub>3</sub> and AGRF<sub>3</sub> models (figure 2), HR was  
235 associated with the greatest increase in mean squared error.

236 **Figure 2 – insert here**

237

## 238 **Discussion**

239 This is the first study to demonstrate that sensor data obtained from commercial  
240 wearable devices (i.e. FC2) can be used to estimate EE with a high degree of validity in a  
241 diverse range of activities. Commercial activity monitors offer a number of benefits over  
242 research-grade devices, including their economic viability, participant acceptance and cloud  
243 storage capabilities [29] and our findings highlight the potential for these inexpensive tools to  
244 more accurately quantify EE. We show that accelerometer data collected from research-grade  
245 devices on the wrist or arm, can be used to predict EE with a high degree of accuracy in a  
246 diverse population.

247 The results of the present study are comparable to the accuracy reported in previous  
248 studies, with overall RMSE reaching 1 MET for the most accurate model (SWRF<sub>2</sub>) and 1.37  
249 METs for the least accurate model (FBRF<sub>2</sub>). Using Actigraph accelerometer data, Ellis et al.  
250 trained random forest models and reported a RMSE of 1 METs for bout predictions from a  
251 hip worn accelerometer, 1.09 from a wrist accelerometer [11]. Staudenmayer et al. developed  
252 artificial neural networks and report a RMSE of 1.22 METs [9]. Montoye et al., showed that  
253 artificial neural networks trained on wrist accelerometer data were more accurate than  
254 corresponding linear models, with RMSE values between 1.26–1.32 METs [12]. Importantly  
255 however, this is the first study to apply machine learning to the sensor data obtained from a  
256 commercially available tracking device (e.g. FC2) in order to improve the accuracy of EE  
257 estimates. In the present study, the RMSE for FBRF models approached the accuracy of the

258 models developed using research grade accelerometers, indicating that models developed  
259 using minute-level data from wearable sensors in combination subject characteristics can be  
260 modelled to accurately predict EE. That said, the activities performed in our protocol are  
261 generally less diverse than some of the aforementioned studies.

262         The models with more accelerometer variables as input features (AGRF and SWRF)  
263 led to the greatest predictive accuracy, indicating the importance of tri-axial accelerometry.  
264 The variable importance plots show that HR was the most important determinant of error  
265 reduction for the SWRF<sub>3</sub> and AGRF<sub>3</sub> models, which reflects the established relationship  
266 between HR and VO<sub>2</sub> [30,31]. In contrast, the FBRF<sub>2</sub> model was less influenced by this  
267 variable and we postulate that this may be related to the sensor used to collect HR estimates.  
268 The polar HR strap, which was used in two AGRF and SWRF models shows near perfect  
269 agreement with electrocardiogram criterion measures [32]. Conversely,  
270 photoplethysmography based HR sensors may produce ‘spurious’ HR measurements [10],  
271 which increases noise in the training and testing data sets and had a detrimental effect on the  
272 predictive accuracy.

273         We normalised HR to each participant’s sitting HR and used this as a predictor in the  
274 developed models. This was motivated by the established relationship between the sitting HR  
275 and the flex point [33]. Given the important predictive role this variable has on the models,  
276 the use of HR above sitting appears to offer a means of capturing some of the individual  
277 variability in the relationship between VO<sub>2</sub> and HR without the need for individual  
278 calibration, as in previous approaches [30]. Despite the importance of HR in achieving the  
279 highest predictive accuracy, we show that AG and SWA models developed without external  
280 HR data can still produce highly accurate estimates of EE, surpassing the manufacturer  
281 estimates of the SWA and FBRF models. Considering the potential burden of additional

282 wearable devices (i.e. chest HR straps), this finding has implications for population research  
283 in which the use of a single device is of considerable appeal.

284         The addition of body composition resulted in a significantly higher MAPE in the case  
285 of SWRF models and this may be explained by the specifics of the random forest algorithm.  
286 Each of the trees grown in the random forest regressor samples approximately 1/3 of the  
287 predictor variables and a bootstrapped sample of the training data, which serves to  
288 decorrelate the trees and limits the likelihood of overfitting. Introducing body composition  
289 could result in a situation in which splits are less likely to include the most relevant predictor  
290 variables such as HR, which theoretically could have a detrimental effect on the global model  
291 [34]. Furthermore, body composition (FFM) is highly correlated with RMR [35], it is likely  
292 that any variance attributable to body composition is already accounted for by normalising  
293 our predictions relative to RMR. Thus, FFM may have a greater predictive ability if absolute  
294 caloric expenditure, rather than METs, was the outcome measure. Regardless, the accurate  
295 and precise measurement of body composition is time consuming, costly and requires  
296 experimental expertise. In this sense, the finding that body composition is not critical to the  
297 predictive accuracy of the random forests has positive implications for the utility of the  
298 models in large studies.

299         We included ambulatory, resting, cycling and housework tasks, which are challenging  
300 to assess using wearable devices owing to the differing accelerometer patterns produced by  
301 each activity [36]. It is notable that the AGRF and SWRF models were statistically equivalent  
302 with the criterion measure (indirect calorimetry) in a many activity modalities. This  
303 demonstrates the potential of a single non-linear model to accurately estimate the  
304 bioenergetic demands of common activity types and overcome the limitations of traditional  
305 linear approaches, which have a tendency to generalise poorly across the spectrum of human  
306 activities [36–38]. However, equivalence was not achieved in all activities; activity specific

307 prediction may be enhanced by generating larger training data or combining activity  
308 classification with regression models.

309         The SWA is considered one of the more valid wearable devices for estimating TDEE  
310 [39–41], and a recent meta-analysis from our group showed that this monitor was one of the  
311 only wrist or arm-worn activity monitoring devices not to systematically over or  
312 underestimate EE overall relative to DLW [4]. Thus, we compared the SWA METs to models  
313 developed in this study to facilitate the interpretation of the developed models. All models  
314 had lower MAPE, RMSE and higher correlations to the criterion measure than the SWA.  
315 These data therefore suggest that machine learning can be applied to sensor data of  
316 commercially available devices to surpass the accuracy of widely used research-grade  
317 devices. The SWA provides sufficiently accurate estimates of EE that it can be used with  
318 measures of body weight/composition to estimate true EI [42]. Unfortunately, studies such as  
319 this are limited in their duration owing to the data storage capabilities and battery life of the  
320 SWA. Our results indicate that it may be possible to replace the SWA with machine learning  
321 models applied to commercial devices, i.e. FC2. Given that these data are accessible  
322 continually from the Fitbit API, this could offer an opportunity for a new generation of  
323 quantitative, long term energy balance research. Mathematical models developed to predict  
324 EI from body weight have been proposed and demonstrate a high degree of accuracy  
325 compared to EI calculated through DLW and DEXA [43,44]. In addition to the cost  
326 associated with these techniques, a recognised limitation of this model is the lack of a  
327 continuous EE estimate. Refining estimates of EE would improve the accuracy of such  
328 models and provide important data on day-to-day variability in physical activity behaviours  
329 and associated EE currently not measurable by the DLW method [45].

330         A benefit of the present study is the expression of EE in METs relative to each  
331 participant's measured RMR. The assumption that 1 MET is equivalent to  $3.5 \text{ ml O}_2/\text{kg}^{-1}$

332  $l/min^{-1}$  can result in substantial bias, depending on the age and body composition of the  
333 subject in question [46]. Secondly, we report the results of two validation methods, a holdout  
334 approach and LOOCV approach. We envisage different experimental protocols in which  
335 participant's data may be available for a calibration procedure prior to beginning an  
336 observation period, as is practiced in the historical 'flex' method [30]. In this situation, the  
337 holdout method may be more reflective of the potential accuracy. It is more probable that  
338 individual calibration would not be possible, and the models would be applied to unseen  
339 participants, in this case a LOOCV is more appropriate. Thirdly, we computed accuracy  
340 statistics from all available minutes in the dataset, rather than aggregating them to 'activity  
341 bouts'. Indeed averaging in this manner has the potential to smooth errors and result in an  
342 artificially low average error. Human activity is performed for different durations and it is  
343 therefore valuable to determine accuracy at the minute-level.

344         Several limitations of this study should be acknowledged, firstly, approximately 70%  
345 of our sample were female. Secondly, the confinement to a laboratory and the utilisation of  
346 only steady state activity minutes may limit the ecological validity of this study and it therefore  
347 remains uncertain how well these models will perform in a less controlled environment with  
348 different activity types. It will be important to validate the models against whole-room  
349 calorimetry and the DLW method. Thirdly, we made no attempt to impute missing data in this  
350 study; devices will be removed or may fail in real life situations this is likely to create missing  
351 data. Considering these limitations, it is of great importance to continuously test and refine the  
352 presented models using data collected from different sedentary and active behaviours,  
353 participants, devices and durations. This is particularly important for commercial devices as  
354 updated and/or new models regularly come to market; nevertheless, the utilisation of three  
355 devices in this study indicates that the modelling approach taken would be applicable to newer  
356 devices.



357

358 **Conclusion**

359           This study demonstrates the potential for machine learning models developed using  
360 minute-level data from wearable sensors in combination subject characteristics to be  
361 modelled to predict EE with minimal bias. Further, machine learning models using outputs  
362 from a commercial activity monitor achieve greater predictive accuracy than the SWA  
363 armband. This methodology opens the possibility for quantitative energy balance research  
364 with affordable, unobtrusive wearable sensors.

365 **References**

- 366 1. Racette SB, Das SK, Bhapkar M, Hadley EC, Roberts SB, Ravussin E, et al. Approaches  
367 for quantifying energy intake and %calorie restriction during calorie restriction interventions  
368 in humans: the multicenter CALERIE study. *AJP Endocrinol Metab* [Internet].  
369 2012;302:E441–8. Available from:  
370 <http://ajpendo.physiology.org/cgi/doi/10.1152/ajpendo.00290.2011>
- 371 2. Black AE, Cole TJ. Within- and between-subject variation in energy expenditure measured  
372 by the doubly-labelled water technique: Implications for validating reported dietary energy  
373 intake. *Eur J Clin Nutr* [Internet]. 2000;54:386–94. Available from:  
374 <http://www.nature.com/doi/10.1038/sj.ejcn.1600970>
- 375 3. Jakicic JM, Marcus M, Gallagher KI, Randall C, Thomas E, Goss FL, et al. Evaluation of  
376 the SenseWear Pro Armband to Assess Energy Expenditure during Exercise. *Med Sci Sport*  
377 *Exerc* [Internet]. United States; 2004;36:897–904. Available from:  
378 [http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00005768-](http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00005768-200405000-00024)  
379 [200405000-00024](http://content.wkhealth.com/linkback/openurl?sid=WKPTLP:landingpage&an=00005768-200405000-00024)
- 380 4. O’Driscoll R, Turicchi J, Beaulieu K, Scott S, Matu J, Deighton K, et al. How well do  
381 activity monitors estimate energy expenditure? A systematic review and meta-analysis of the  
382 validity of current technologies. *Br J Sports Med* [Internet]. BMJ Publishing Group Ltd and  
383 British Association of Sport and Exercise Medicine; 2018 [cited 2018 Oct 1];77:bjsports-  
384 2018-099643. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/30194221>
- 385 5. Freedson PS, Melanson E, Sirard J. Calibration of the Computer Science and Applications,  
386 Inc. accelerometer. *Med Sci Sports Exerc*. United States; 1998;30:777–81.
- 387 6. Crouter SE, Bassett DRJ. A new 2-regression model for the Actical accelerometer. *Br J*  
388 *Sports Med*. England; 2008;42:217–24.
- 389 7. Bonomi AG, Plasqui G, Goris AHC, Westerterp KR. Improving assessment of daily

- 390 energy expenditure by identifying types of physical activity with a single accelerometer. *J*  
391 *Appl Physiol* [Internet]. 2009;107:655–61. Available from:  
392 <http://www.physiology.org/doi/10.1152/jappphysiol.00150.2009>
- 393 8. Hills AP, Mokhtar N, Byrne NM. Assessment of Physical Activity and Energy  
394 Expenditure: An Overview of Objective Measures. *Front Nutr* [Internet]. 2014;1:1–16.  
395 Available from: <http://journal.frontiersin.org/article/10.3389/fnut.2014.00005/abstract>
- 396 9. Staudenmayer J, Pober D, Crouter S, Bassett D, Freedson P. An artificial neural network to  
397 estimate physical activity energy expenditure and identify physical activity type from an  
398 accelerometer. *J Appl Physiol* [Internet]. United States: American Physiological Society;  
399 2009 [cited 2019 Feb 5];107:1300–7. Available from:  
400 <http://www.ncbi.nlm.nih.gov/pubmed/19644028>
- 401 10. Reddy RK, Pooni R, Zaharieva DP, Senf B, El Youssef J, Dassau E, et al. Accuracy of  
402 Wrist-Worn Activity Monitors During Common Daily Physical Activities and Types of  
403 Structured Exercise: Evaluation Study. *JMIR mHealth uHealth* [Internet]. *JMIR mHealth and*  
404 *uHealth*; 2018 [cited 2018 Dec 16];6:e10338. Available from:  
405 <https://mhealth.jmir.org/2018/12/e10338/>
- 406 11. Ellis K, Kerr J, Godbole S, Lanckriet G, Wing D, Marshall S. A random forest classifier  
407 for the prediction of energy expenditure and type of physical activity from wrist and hip  
408 accelerometers. *Physiol Meas* [Internet]. England; 2014 [cited 2019 Feb 5];35:2191–203.  
409 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/25340969>
- 410 12. Montoye AHKK, Begum M, Henning Z, Pfeiffer KA. Comparison of linear and non-  
411 linear models for predicting energy expenditure from raw accelerometer data. *Physiol Meas*.  
412 IOP Publishing; 2017;38:343.
- 413 13. ACSM. Exercise Preparticipation Health Screen Recommendations. 2018 [cited 2018 Feb  
414 20]; Available from: <http://www.acsm.org/docs/default-source/publications/acsm-101->

- 415 prescreeninginfographiccolorlegal-2015-12-15-v02.pdf?sfvrsn=2
- 416 14. William E. Siri. BODY COMPOSITION FROM FLUID SPACES AND DENSITY: ,  
417 ANALYSIS OF METHODS. *Adv Biol Med Phy.* 1956;
- 418 15. Benedetto S, Caldato C, Bazzan E, Greenwood DC, Pensabene V, Actis P. Assessment of  
419 the fitbit charge 2 for monitoring heart rate. *PLoS One [Internet].* 2018;13:e0192691.  
420 Available from: <http://dx.plos.org/10.1371/journal.pone.0192691>
- 421 16. Sanchez-Delgado G, Alcantara JMA, Ortiz-Alvarez L, Xu H, Martinez-Tellez B, Labayen  
422 I, et al. Reliability of resting metabolic rate measurements in young adults: Impact of  
423 methods for data analysis. *Clin Nutr [Internet].* 2018 [cited 2019 Apr 24];37:1618–24.  
424 Available from: <https://linkinghub.elsevier.com/retrieve/pii/S0261561417302674>
- 425 17. Müller MJ, Bosy-Westphal A, Klaus S, Kreymann G, Lührmann PM, Neuhäuser-  
426 Berthold M, et al. World Health Organization equations have shortcomings for predicting  
427 resting energy expenditure in persons from a modern, affluent population: generation of a  
428 new reference standard from a retrospective analysis of a German database of resting energy  
429 expe. *Am J Clin Nutr [Internet].* 2004;80:1379–90. Available from:  
430 <http://www.ncbi.nlm.nih.gov/pubmed/15531690>
- 431 18. Perez-Suarez I, Martin-Rincon M, Gonzalez-Henriquez JJ, Fezzardi C, Perez-Regalado S,  
432 Galvan-Alvarez V, et al. Accuracy and Precision of the COSMED K5 Portable Analyser.  
433 *Front Physiol [Internet]. Frontiers;* 2018 [cited 2019 Apr 1];9:1764. Available from:  
434 <https://www.frontiersin.org/article/10.3389/fphys.2018.01764/full>
- 435 19. Groepenhoff H, de Jeu RC, Schot R. Vyntus CPX compared to Oxycon pro shows equal  
436 gas-exchange and ventilation during exercise. *Respir Funct Technol [Internet]. European  
437 Respiratory Society;* 2017 [cited 2019 Apr 24]. p. PA3002. Available from:  
438 <http://erj.ersjournals.com/lookup/doi/10.1183/1393003.congress-2017.PA3002>
- 439 20. Péronnet F, Massicotte D. Table of nonprotein respiratory quotient: an update. *Can J*

- 440 Sport Sci [Internet]. 1991 [cited 2019 Jun 11];16:23–9. Available from:  
441 <http://www.ncbi.nlm.nih.gov/pubmed/1645211>
- 442 21. Staudenmayer J, Pober D, Crouter S, Bassett D, Freedson P. An artificial neural network  
443 to estimate physical activity energy expenditure and identify physical activity type from an  
444 accelerometer. *J Appl Physiol*. 2009;107:1300–7.
- 445 22. Breiman L. Random Forests. *Mach Learn [Internet]*. Kluwer Academic Publishers; 2001  
446 [cited 2019 Feb 9];45:5–32. Available from:  
447 <http://link.springer.com/10.1023/A:1010933404324>
- 448 23. Tibshirani TH and R. *Introduction to Statistical Learning*. Japan Med. Assoc. J. 2011.
- 449 24. Liaw A, Wiener M, Andy Liaw M. Package: “randomForest”. 2018 [cited 2019 May 13];  
450 Available from: <https://www.stat.berkeley.edu/~breiman/RandomForests/>
- 451 25. Bai Y, Hibbing P, Mantis C, Welk GJ. Comparative evaluation of heart rate-based  
452 monitors: Apple Watch vs Fitbit Charge HR. *J Sports Sci [Internet]*. Routledge;  
453 2018;36:1734–41. Available from: <https://doi.org/10.1080/02640414.2017.1412235>
- 454 26. Hamner B, Frasco M, Ledell E. Package “Metrics” Title Evaluation Metrics for Machine  
455 Learning [Internet]. 2018. Available from: [https://cran.r-](https://cran.r-project.org/web/packages/Metrics/Metrics.pdf)  
456 [project.org/web/packages/Metrics/Metrics.pdf](https://cran.r-project.org/web/packages/Metrics/Metrics.pdf)
- 457 27. Lee J-MM, Kim Y-WY, Welk GJ. Validity of consumer-based physical activity monitors.  
458 *Med Sci Sports Exerc [Internet]*. United States, United States; 2014 [cited 2017 Nov  
459 20];46:1840–8. Available from:  
460 <http://search.ebscohost.com/login.aspx?direct=true&db=sph&AN=96789748&site=ehost-live>
- 461 28. Team R. Core. R: A language and environment for statistical computing. 2018;0:201.
- 462 29. Wright SP, Hall Brown TS, Collier SR, Sandberg K. How consumer physical activity  
463 monitors could transform human physiology research. *Am J Physiol Integr Comp Physiol*  
464 [Internet]. S.P. Wright, Georgetown Univ., 232 Building D, 4000 Reservoir Rd., NW,

- 465 Washington, DC 20057, United States. E-mail: spw44@georgetown.edu, United States:  
466 American Physiological Society (E-mail: subscrip@the-aps.org); 2017 [cited 2017 Nov  
467 17];312:R358–67. Available from:  
468 <http://ajpregu.physiology.org/lookup/doi/10.1152/ajpregu.00349.2016>
- 469 30. Ceesay SM, Prentice AM, Day KC, Murgatroyd PR, Goldberg GR, Scott W, et al. The  
470 use of heart rate monitoring in the estimation of energy expenditure : a validation study using  
471 indirect whole-body calorimetry. *British J Nutr* [Internet]. England; 1989 [cited 2017 Oct  
472 23];61:175–86. Available from:  
473 [http://www.journals.cambridge.org/abstract\\_S0007114589000267](http://www.journals.cambridge.org/abstract_S0007114589000267)
- 474 31. Leonard WR. Measuring human energy expenditure: What have we learned from the  
475 flex-heart rate method? *Am J Hum Biol*. 2003;15:479–89.
- 476 32. Gillinov S, ETIWY M, Wang R, Blackburn G, PHELAN D, Gillinov AM, et al. Variable  
477 Accuracy of Wearable Heart Rate Monitors during Aerobic Exercise. *Med Sci Sport Exerc*  
478 [Internet]. Lippincott Williams and Wilkins; 2017 [cited 2017 Nov 17];49:1697–703.  
479 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/28709155>
- 480 33. Rennie KL, Hennings SJ, Mitchell J, Wareham NJ. Estimating energy expenditure by  
481 heart-rate monitoring without individual calibration. *Med Sci Sports Exerc* [Internet].  
482 2001;33:939–45. Available from:  
483 [http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=11404659&retmo  
484 de=ref&cmd=prlinks%5Cnpapers2://publication/uuid/5CD2363C-779F-4417-B26D-  
485 4C2A80457601](http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=11404659&retmode=ref&cmd=prlinks%5Cnpapers2://publication/uuid/5CD2363C-779F-4417-B26D-4C2A80457601)
- 486 34. Hastie T, Tibshirani R, Friedman J. *The Elements of Statistical Learning The Elements of*  
487 *Statistical Learning*. 2017; Available from: <https://web.stanford.edu/~hastie/Papers/ESLII.pdf>
- 488 35. Hopkins M, Finlayson G, Duarte C, Whybrow S, Ritz P, Horgan GW, et al. Modelling  
489 the associations between fat-free mass, resting metabolic rate and energy intake in the context

- 490 of total energy balance. *Int J Obes* [Internet]. Nature Publishing Group; 2016;40:312–8.  
491 Available from: <http://dx.doi.org/10.1038/ijo.2015.155>
- 492 36. Chowdhury EA, Western MJ, Nightingale TE, Peacock OJ, Thompson D. Assessment of  
493 laboratory and daily energy expenditure estimates from consumer multisensor physical  
494 activity monitors. *PLoS One* [Internet]. United States, United States: Public Library of  
495 Science (E-mail: [plos@plos.org](mailto:plos@plos.org)); 2017 [cited 2017 Nov 9];12:e0171720. Available from:  
496 <http://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0171720&type=printable>  
497 e
- 498 37. Lee J-M, Kim Y-W, Welk GJ. Validity and Utility of Consumer-Based Physical Activity.  
499 *ACSM's Heal Fit J* [Internet]. 2014;18:16–21. Available from:  
500 <http://search.ebscohost.com/login.aspx?direct=true&db=sph&AN=96789748&site=ehost-live>
- 501 38. Wahl Y, Düking P, Droszez A, Wahl P, Mester J, Y. W, et al. Criterion-validity of  
502 commercially available physical activity tracker to estimate step count, covered distance and  
503 energy expenditure during sports conditions. *Front Physiol* [Internet]. Y. Wahl, Institute of  
504 Biomechanics and Orthopedics, German Sport University Cologne, Cologne, Germany. E-  
505 mail: [y.wahl@dshs-koeln.de](mailto:y.wahl@dshs-koeln.de), Switzerland: Frontiers Media S.A. (E-mail:  
506 [info@frontiersin.org](mailto:info@frontiersin.org)); 2017 [cited 2017 Nov 13];8:725. Available from:  
507 <http://journal.frontiersin.org/article/10.3389/fphys.2017.00725/full>
- 508 39. Slinde F, Bertz F, Winkvist A, Ellegård L, Olausson H, Brekke HK. Energy expenditure  
509 by multisensor armband in overweight and obese lactating women validated by doubly  
510 labeled water. *Obesity* [Internet]. United States; 2013 [cited 2017 Dec 10];21:2231–5.  
511 Available from: <http://doi.wiley.com/10.1002/oby.20363>
- 512 40. Casiraghi F, Lertwattanak R, Luzi L, Chavez AO, Davalli AM, Naegelin T, et al.  
513 Energy Expenditure Evaluation in Humans and Non-Human Primates by SenseWear  
514 Armband. Validation of Energy Expenditure Evaluation by SenseWear Armband by Direct

- 515 Comparison with Indirect Calorimetry. PLoS One [Internet]. F. Folli, Department of  
516 Medicine/Division of Diabetes, University of Texas Health Science Center at San Antonio,  
517 San Antonio, TX, United States. E-mail: folli@uthscsa.edu, United States: Public Library of  
518 Science (185 Berry Street, Suite 1300, San Francisco CA 94107, United States);  
519 2013;8:e73651. Available from: <http://www.plosone.org/article/fetchObject.action>
- 520 41. Johannsen DL, Calabro MA, Stewart J, Franke W, Rood JC, Welk GJ. Accuracy of  
521 armband monitors for measuring daily energy expenditure in healthy adults. *Med Sci Sports*  
522 *Exerc* [Internet]. G. J. Welk, 257 Forker Bldg., Ames, IA 50011, United States. E-mail:  
523 gwelk@iastate.edu: Lippincott Williams and Wilkins (530 Walnut Street, P O Box 327,  
524 Philadelphia PA 19106-3621, United States); 2010;42:2134–40. Available from:  
525 <http://search.ebscohost.com/login.aspx?direct=true&db=sph&AN=55265135&site=ehost-live>
- 526 42. Shook RP, Hand GA, O'Connor DP, Thomas DM, Hurley TG, Hébert JR, et al. Energy  
527 Intake Derived from an Energy Balance Equation, Validated Activity Monitors, and Dual X-  
528 Ray Absorptiometry Can Provide Acceptable Caloric Intake Data among Young Adults. *J*  
529 *Nutr* [Internet]. United States; 2018;148:490–6. Available from:  
530 <https://academic.oup.com/jn/article/148/3/490/4930797>
- 531 43. Sanghvi A, Redman LM, Martin CK, Ravussin E, Hall KD. Validation of an inexpensive  
532 and accurate mathematical method to measure long-term changes in free-living energy  
533 intake. *Am J Clin Nutr*. United States; 2015;102:353–8.
- 534 44. Thomas DDM, Schoeller DA, Redman LA, Martin CK, Levine JA, Heymsfield SB. A  
535 computational model to determine energy intake during weight loss. *Am J ...* [Internet].  
536 United States; 2010;92:1326–31. Available from:  
537 <http://ajcn.nutrition.org/content/92/6/1326.short>
- 538 45. Schoeller DA, Ravussin E, Schutz Y, Acheson KJ, Baertschi P, Jequier E, et al. Energy  
539 expenditure by doubly labeled water: validation in humans and proposed calculation. *Am J*



540 Physiol - Regul Integr Comp Physiol [Internet]. 1986;250:R823–30. Available from:  
541 <http://ajpregu.physiology.org/content/250/5/R823.abstract>  
542 46. Byrne NM, Hills AP, Hunter GR, Weinsier RL, Schutz Y. Metabolic equivalent: one size  
543 does not fit all. J Appl Physiol [Internet]. 2005 [cited 2019 May 7];99:1112–9. Available  
544 from: <http://www.physiology.org/doi/10.1152/jappphysiol.00023.2004>  
545

546 **Tables**

547 Table 1: Predictive variables included in each of the random forest models.

<b>Fitbit Charge</b> <b>2</b>		
	<b>Model</b>	
	FBRF <sub>1</sub>	<p><b>Device outputs:</b></p> <p>Fitbit HR above sitting HR, Steps,</p> <p><b>Subject characteristics:</b></p> <p>Age, Gender, Height, weight</p>
	FBRF <sub>2</sub>	1 + FM (kg), FFM (kg)
<b>Sensewear</b> <b>Armband Mini</b>		
	SWRF <sub>1</sub>	<p><b>Device outputs:</b></p> <p>Average (Axis: X, Y, Z), Peaks (X, Y, Z), Mean absolute deviation (X, Y, Z), steps/min-1, Near body temperature average, skin temperature average, Galvanic skin response.</p> <p><b>Subject characteristics:</b></p> <p>age, gender, Height, weight,</p>
	SWRF <sub>2</sub>	1 + polar HR above sitting HR
	SWRF <sub>3</sub>	2 + FM(kg), FFM (kg)
Actigraph GT3- x	AGRF <sub>1</sub>	<b>Time domain, multi-axis (X, Y, Z) and first order differential (XYZ) features:</b>

		<p>minimum, maximum, mean, standard deviation, correlation (XY, XZ, YZ), Median 0 crossings, percentiles (10, 25, 50, 75, 90<sup>th</sup>)</p> <p><b>Frequency domain multi-axis (X, Y, Z) and first order differential (XYZ) features:</b></p> <p>dominant frequency, dominant frequency magnitude</p> <p><b>Subject characteristics:</b></p> <p>Gender, age, height, weight,</p>
	AGRF <sub>2</sub>	<b>1 + Polar Heart rate above sitting HR</b>
	AGRF <sub>3</sub>	2 + FM (kg), FFM (kg)

548

549 Abbreviations: Fitbit random forest (FBRF), Sensewear random forest (SWRF), Actigraph  
550 random forest (AGRF).

551 Table 2: Accuracy statistics computed using the LOOCV validation approach. Criterion  
 552 METs Refers to METs calculated from indirect calorimetry, Predicted METs refers to model  
 553 prediction. Minutes pooled refers to the number of minutes used for validation and  
 554 participants refers to the number of participants included in each validation. RMSE, MAPE  
 555 and correlation are presented with 95% confidence intervals. Equivalence refers to the results  
 556 of the equivalence tests.  
 557

	Minutes	Participants	Criterion METs	Predicted METs	RMSE (METs)	MAPE (%)	Correlation (r)	Equivalence
SWA	2188	59	4.19 ± 2.61	3.72 ± 2.40	1.8	33.57	0.76	not equivalent
AGR F <sub>1</sub>	2161	58	4.20 ± 2.61	4.18 ± 2.32	1.18	20.9	0.89	equivalent
AGR F <sub>2</sub>	2125	57	4.22 ± 2.63	4.21 ± 2.37	1.03	18.31	0.92	equivalent
AGR F <sub>3</sub>	2049	55	4.22 ± 2.62	4.21 ± 2.37	1.02	18.52	0.92	equivalent
FBRF <sub>1</sub>	2077	57	4.19 ± 2.63	4.16 ± 2.30	1.36	28.74	0.86	equivalent
FBRF <sub>2</sub>	2001	55	4.19 ± 2.63	4.14 ± 2.20	1.37	30.59	0.85	equivalent
SWR F <sub>1</sub>	2188	59	4.19 ± 2.61	4.19 ± 2.24	1.24	23.61	0.88	equivalent
SWR F <sub>2</sub>	2153	58	4.21 ± 2.62	4.22 ± 2.35	1	18.82	0.92	equivalent
SWR F <sub>3</sub>	2077	56	4.21 ± 2.62	4.22 ± 2.34	1.02	19.38	0.92	equivalent

558

559

560 Abbreviations: Metabolic equivalents (METs), Root mean squared error (RMSE), Mean

561 absolute percentage error (MAPE) , Actigraph random forest (AGRF), Fitbit random forest

562 (FBRF), Sensewear random forest (SWRF).

563



	SWA	236	59	5.24 ± 0.72	4.15 ± 0.85	1.65	26.14	-0.26	not equivalent
	AGRF 1	232	58	5.23 ± 0.72	4.68 ± 0.58	1.03	14.14	0.11	not equivalent
	AGRF 2	228	57	5.23 ± 0.73	4.82 ± 0.65	0.96	13.27	0.21	equivalent
	AGRF 3	220	55	5.26 ± 0.72	4.85 ± 0.66	0.96	13.01	0.21	equivalent
	FBRF1	216	57	5.19 ± 0.72	4.36 ± 0.46	1.16	17.01	0.1	not equivalent
	FBRF2	208	55	5.22 ± 0.72	4.33 ± 0.49	1.2	18.87	0.14	not equivalent
	SWRF 1	236	59	5.24 ± 0.73	4.92 ± 0.97	1.28	18.36	-0.05	equivalent
	SWRF 2	232	58	5.24 ± 0.73	4.97 ± 0.65	0.94	13.86	0.14	equivalent
	SWRF 3	224	56	5.27 ± 0.72	5 ± 0.65	0.93	13.48	0.14	equivalent
Jog									
	SWA	195	49	8.57 ± 1.21	8.12 ± 1.52	2.01	19.3	-0.01	equivalent
	AGRF 1	195	49	8.57 ± 1.21	8.73 ± 1.12	1.54	15.87	0.14	equivalent
	AGRF 2	195	49	8.57 ± 1.21	8.63 ± 1.15	1.38	13.81	0.32	equivalent
	AGRF 3	187	47	8.54 ± 1.23	8.59 ± 1.14	1.33	13.56	0.36	equivalent
	FBRF1	191	48	8.55 ± 1.22	8.68 ± 1.19	1.5	15.3	0.22	equivalent
	FBRF2	183	46	8.52 ± 1.24	8.46 ± 1.18	1.54	15.57	0.18	equivalent
	SWRF 1	195	49	8.57 ± 1.21	8.48 ± 1.15	1.34	12.33	0.36	equivalent
	SWRF 2	195	49	8.57 ± 1.21	8.59 ± 1.08	1.15	11.26	0.5	equivalent
	SWRF 3	187	47	8.54 ± 1.23	8.5 ± 1.1	1.17	11.45	0.5	equivalent
Jog incline									
	SWA	120	30	10.07 ± 1.32	8.04 ± 1.75	3.03	24.69	-0.06	not equivalent
	AGRF 1	120	30	10.07 ± 1.32	8.59 ± 1.43	2.42	19.33	0.03	not equivalent
	AGRF 2	120	30	10.07 ± 1.32	8.86 ± 1.45	2.19	18.03	0.13	not equivalent
	AGRF 3	116	29	10.06 ± 1.34	8.98 ± 1.37	2.04	17.12	0.17	not equivalent
	FBRF1	120	30	10.07 ± 1.32	8.86 ± 1.54	2.1	15.9	0.28	not equivalent
	FBRF2	116	29	10.06 ± 1.34	8.7 ± 1.5	2.21	17.4	0.24	not equivalent
	SWRF 1	120	30	10.07 ± 1.32	8.36 ± 1.48	2.45	19.15	0.21	not equivalent
	SWRF 2	120	30	10.07 ± 1.32	8.88 ± 1.24	1.83	15.2	0.41	not equivalent
	SWRF 3	116	29	10.06 ± 1.34	8.9 ± 1.26	1.81	15.03	0.42	not equivalent
Cycle low									
	SWA	233	59	4.15 ± 1.08	2.53 ± 0.96	1.93	40.49	0.46	not equivalent
	AGRF 1	229	58	4.13 ± 1.07	4.33 ± 0.74	1.25	25.11	0.11	equivalent
	AGRF 2	225	57	4.14 ± 1.07	4.33 ± 0.76	1.03	20.8	0.42	equivalent
	AGRF 3	217	55	4.14 ± 1.08	4.32 ± 0.75	1.06	21.6	0.39	equivalent
	FBRF1	220	56	4.11 ± 1.09	3.77 ± 1.08	1.54	29.31	0.03	not equivalent
	FBRF2	212	54	4.1 ± 1.1	3.63 ± 0.92	1.45	26.67	0.09	not equivalent
	SWRF 1	233	59	4.15 ± 1.08	4.25 ± 0.75	1.11	22.16	0.3	equivalent

	SWRF 2	229	58	4.17 ± 1.08	4.23 ± 0.86	0.93	17.81	0.56	equivalent
	SWRF 3	221	56	4.17 ± 1.09	4.2 ± 0.86	0.97	18.39	0.53	equivalent
Cycle mid									
	SWA	225	58	5.21 ± 1.43	3.31 ± 1.75	2.55	41.99	0.44	not equivalent
	AGRF 1	225	58	5.21 ± 1.43	4.53 ± 0.65	1.61	20.11	0.18	not equivalent
	AGRF 2	221	57	5.22 ± 1.44	4.78 ± 0.81	1.32	16.67	0.5	not equivalent
	AGRF 3	213	55	5.25 ±1.46	4.82 ± 0.83	1.34	17.08	0.5	not equivalent
	FBRF1	207	55	5.19 ± 1.44	3.71 ± 1.13	2.27	30.96	0.13	not equivalent
	FBRF2	199	53	5.22 ± 1.47	3.61 ± 0.94	2.27	30.75	0.16	not equivalent
	SWRF 1	225	58	5.21 ± 1.43	4.16 ± 0.84	1.79	22.21	0.26	not equivalent
	SWRF 2	221	57	5.22 ± 1.44	4.59 ± 1.03	1.47	18.89	0.47	not equivalent
	SWRF 3	213	55	5.25 ± 1.46	4.58 ± 1.01	1.49	19.03	0.46	not equivalent
Folding									
	SWA	236	59	2.75 ± 0.6	4.29 ± 1.77	2.37	72.37	0.12	not equivalent
	AGRF 1	232	58	2.75 ± 0.6	2.91 ± 0.35	0.61	18.98	0.32	equivalent
	AGRF 2	228	57	2.75 ± 0.61	2.95 ± 0.38	0.59	18.8	0.44	equivalent
	AGRF 3	220	55	2.75 ± 0.61	2.96 ± 0.39	0.63	19.75	0.36	equivalent
	FBRF1	228	57	2.74 ± 0.61	3.5 ± 0.71	1.19	37.21	0.05	not equivalent
	FBRF2	220	55	2.75 ± 0.61	3.58 ± 0.71	1.26	39.8	-0.01	not equivalent
	SWRF 1	236	59	2.75 ± 0.6	3.43 ± 0.35	1.04	33.41	0.08	not equivalent
	SWRF 2	232	58	2.74 ± 0.6	3.37 ± 0.69	0.98	30.06	0.32	not equivalent
	SWRF 3	224	56	2.75 ± 0.61	3.42 ± 0.7	1.03	32.08	0.29	not equivalent
Sweepin g									
	SWA	235	59	3.12 ± 0.71	3.52 ± 1.47	1.59	41.49	0.13	not equivalent
	AGRF 1	232	58	3.13 ± 0.71	3.64 ± 0.68	1.03	29.33	0.17	not equivalent
	AGRF 2	226	57	3.12 ± 0.71	3.57 ± 0.68	0.95	27.09	0.28	not equivalent
	AGRF 3	218	55	3.13 ± 0.71	3.56 ± 0.63	0.93	27	0.24	not equivalent
	FBRF1	228	57	3.13 ± 0.72	3.95 ± 0.67	1.2	35.29	0.21	not equivalent
	FBRF2	220	55	3.13 ± 0.71	3.96 ± 0.69	1.22	35.33	0.2	not equivalent
	SWRF 1	235	59	3.13 ± 0.72	3.65 ± 0.83	1.14	30.04	0.15	not equivalent
	SWRF 2	230	58	3.11 ± 0.71	3.51 ± 0.92	0.97	24.11	0.43	not equivalent
	SWRF 3	222	56	3.12 ± 0.71	3.57 ± 0.95	1.02	25.41	0.42	not equivalent

571  
572  
573

574 Abbreviations: Metabolic equivalents (METs), Root mean squared error (RMSE), Mean  
575 absolute percentage error (MAPE), Fitbit random forest (FBRF), Sensewear random forest  
576 (SWRF), Actigraph random forest (AGRF).

577



## 578 Legends for figures

579 Figure 1. A scatter plot of the Measured METs (Vyntus CPX) and predicted METs. Diagonal  
580 lines represent lines of identity. Data are shown for the Actigraph random forest model  
581 (AGRF<sub>2</sub>, top left), Fitbit random forest model (FBRF<sub>1</sub>, top right), Sensewear armband (SWA,  
582 bottom left) and Sensewear armband random forest (SWRF<sub>2</sub>, bottom right).

583

584 Figure 2. Variable importance plots detailing the increase in mean squared error associated  
585 with permutation (random shuffling) of a single variable. Data are shown in order of  
586 importance for the first 20 variables. Data are shown for models: FBRF<sub>2</sub> (blue dots, left),  
587 SWRF<sub>3</sub> (red dots, middle), and AGRF<sub>3</sub> (green dots, right).

588 Abbreviations: HR = Heart rate, FM (kg)= Fat mass (kg), FFM (KG) = Fat free mass (kg),

589 MAD = Mean absolute deviation, FOD = First order differential.



