

**Comparing two wrist-worn accelerometers (Axivity AX3 and Matrix 003) for measuring  
movement behaviours in British and Chinese older adults**

**Date of submission:** Monday 2<sup>nd</sup> June 2025

## 1 **Abstract**

2 **Introduction:** Two nationally representative cohorts, the English Longitudinal Study of Ageing  
3 (ELSA) and the China Health and Retirement Longitudinal Study (CHARLS), recently  
4 introduced accelerometry to measure movement behaviours. However, the use of different  
5 brands (Axivity AX3 and Matrix 003) may hinder data harmonisation. This study assessed  
6 whether the raw acceleration data and machine learning-derived physical activity and sleep  
7 outcomes were equivalent between these two accelerometers in both British and Chinese adults.

8 **Methods:** Eighty-five British and 117 Chinese adults aged  $\geq 50$  years wore both accelerometers  
9 in a random order on their dominant wrist for up to eight days. Data were processed using UK-  
10 derived machine learning algorithms to generate outcomes such as average acceleration (mg),  
11 time in 24-hour movement behaviours (hours/day), daily step count, and sleep duration  
12 (hours/night). Equivalency was assessed using 95% equivalence tests ( $\pm 10\%$  equivalence zone).

13 **Results:** In both British and Chinese adults, average acceleration, sedentary time, time in bed,  
14 and sleep duration were equivalent between the two accelerometers, while time in moderate-  
15 vigorous physical activity (MVPA) was not ( $\pm 17.7\%$  in British and  $\pm 28.2\%$  in Chinese adults).  
16 Time in light physical activity (LPA) was equivalent in British ( $\pm 6.2\%$ ) but borderline in  
17 Chinese adults ( $\pm 10.3\%$ ), whereas the opposite was observed for daily step count ( $\pm 10.5\%$  in  
18 British and  $\pm 2.9\%$  in Chinese adults).

19 **Conclusion:** Average acceleration was comparable between the Axivity AX3 and the Matrix 003  
20 in both British and Chinese adults. Machine learning-derived outcomes were also largely  
21 comparable; however, the cross-nationality differences highlight the need for further population-  
22 specific algorithm development and validation.

23 **Word count:** 250 words

24 **Keywords:** raw acceleration; machine learning; physical activity; sleep; data harmonisation;  
25 cross-nationality comparison.

26 **Highlights**

27 1) The harmonisation of accelerometry data across large-scale epidemiological cohorts is  
28 currently being hindered by different brands and processing methods.

29 2) Average acceleration – a widely used measure of overall physical activity – was  
30 comparable between the Axivity AX3 and the Matrix 003 in both British and Chinese  
31 adults aged 50 years or older; this suggests that data from these two different brands of  
32 wrist-worn accelerometer can be pooled and/or compared.

33 3) Machine learning-derived physical activity and sleep outcomes were also largely  
34 comparable; however, the cross-nationality differences observed highlight the need for  
35 further population-specific algorithm development and validation.

## 36 **Introduction**

37 Less time spent being physically active, more time spent being sedentary, and suboptimal sleep  
38 duration (outside 6-8 hours per night) are associated with premature mortality and a higher risk  
39 of many chronic diseases (Bull et al., 2020; Chaput et al., 2020). However, the majority of  
40 evidence to date is based on self-reported measures of these movement behaviours. Although  
41 self-report has been sufficient to demonstrate the importance of movement behaviours for health,  
42 and can provide valuable contextual information, it lacks the precision to provide specific  
43 quantitative recommendations about the optimal amounts and/or composition of movement  
44 behaviours for health. For example, studies using accelerometry to examine the association  
45 between physical activity and premature mortality have reported effect sizes almost double the  
46 size of studies using self-report (Ekelund et al., 2019; Wasfy & Lee, 2022).

47 An increasing number of large-scale epidemiological cohorts are implementing wrist-worn  
48 accelerometry to measure movement behaviours. However, differences in accelerometer brand  
49 and processing method create challenges for data harmonisation. The UK Biobank (Doherty et  
50 al., 2017), the China Kadoorie Biobank (Chen et al., 2023), and more recently, the English  
51 Longitudinal Study of Ageing (ELSA; 2021-23) (Steptoe et al., 2013) all used the Axivity AX3  
52 wrist-worn accelerometer. In contrast, the US National Health and Nutrition Examination Survey  
53 (NHANES; 2011-14) used the ActiGraph GT3X+, the British Whitehall II Study (Menai et al.,  
54 2017) used the GENEActiv, and the China Health and Retirement Longitudinal Study  
55 (CHARLS; 2021-23) (Zhao et al., 2014) developed their own wrist-worn device, the Matrix 003,  
56 in collaboration with a Chinese manufacturer. Using open-source software (van Hees et al.,  
57 2013; 2014), it is possible to process the raw acceleration data from different accelerometer  
58 brands and directly compare them (Rowlands et al., 2019). In addition, there is now the

59 availability of open-source machine learning algorithms to estimate the following: 1) time spent  
60 in 24-hour movement behaviours – i.e., moderate-vigorous physical activity (MVPA), light  
61 physical activity (LPA), sedentary behaviour (SB), and time in bed (Walmsley et al., 2021); 2)  
62 sleep duration and efficiency (Yuan et al., 2024); and 3) step count and cadence (Small et al.,  
63 2024). However, it remains to be seen whether these machine learning algorithms can be applied  
64 to different accelerometer brands in different populations.

65 This study aimed to test the first two components of the V3 framework for evaluating Biometric  
66 Monitoring Technologies (BioMeTs) such as accelerometers (Goldsack et al., 2020). The  
67 primary aim (verification) was to establish whether the raw acceleration data – i.e., the volume  
68 and intensity distribution of physical activity – from two different brands of accelerometer (the  
69 well-established Axivity AX3 and the newly developed Matrix 003) worn on the dominant wrist  
70 can be considered equivalent in both British and Chinese older adults (aged 50 years or older).  
71 The secondary aim (analytical validation) was to establish whether the machine learning-derived  
72 physical activity and sleep outcomes can also be considered equivalent in both British and  
73 Chinese older adults.

## 74 **Methods**

### 75 *Free-living validation study in British adults*

76 A convenience sample of 85 ambulant adults aged 50 years or older living in the UK was  
77 recruited by email and word of mouth. All participants provided written informed consent, and  
78 the study was approved by the [REDACTED] Research Ethics Committee.  
79 Data were collected between October 2021 and November 2022.

80 Participants self-reported their height and weight, to calculate body mass index (BMI) in kg/m<sup>2</sup>,  
81 and any long-standing (≥3 months) mobility limitations (Supplementary Table S1). Participants  
82 were asked to start wearing both accelerometers immediately after receiving them in the post and  
83 wear them on their dominant wrist 24 hours per day for eight consecutive days. The positional  
84 order of the two accelerometers on the wrist was randomised between participants  
85 (Supplementary Figure S1). After eight days, the █████ Courier Service collected the  
86 accelerometers. Participants were informed that they could wear the waterproof accelerometers  
87 when bathing or swimming but not in extremely high temperature or pressure environments (e.g.,  
88 in a sauna or when diving). Participants were asked to carry on with their normal activities whilst  
89 wearing the accelerometers and did not receive feedback on their activity levels until after they  
90 were returned.

#### 91 *Free-living validation study in Chinese adults*

92 A convenience sample of 117 ambulant adults aged 50 years or older living in China was  
93 recruited by visiting two urban and two rural communities in or near █████. All participants  
94 provided written informed consent, and the study was approved by the █████ Ethics  
95 Review Committee. Data were collected between May and November 2023.

96 Participants self-reported their height and weight, to calculate BMI in kg/m<sup>2</sup>, and any long-  
97 standing (≥3 months) mobility limitations (Supplementary Table S1). The researchers placed  
98 both accelerometers on the participants' dominant wrist in a randomised positional order  
99 (Supplementary Figure S1). Participants were asked to wear the accelerometers 24 hours per day  
100 for two nights and one day. After the second night, the researchers collected the accelerometers.  
101 Participants were informed that they could wear the waterproof accelerometers when bathing or  
102 swimming but not in extremely high temperature or pressure environments (e.g., in a sauna or

103 when diving). Participants were asked to carry on with their normal activities whilst wearing the  
104 accelerometers and did not receive feedback on their activity levels until after they were  
105 collected.

#### 106 *The Axivity AX3*

107 The Axivity (Axivity Ltd, ██████████ UK) is a wrist-worn triaxial accelerometer that has been  
108 used to measure movement behaviours in large-scale epidemiological cohorts such as the UK  
109 Biobank (Doherty et al., 2017), the China Kadoorie Biobank (Chen et al., 2023), and the tenth  
110 wave of ELSA (Steptoe et al., 2013). In British adults, the Axivity was set to start recording at  
111 10am two working days after postal dispatch and stop recording seven full days later, whereas in  
112 Chinese adults it was set to start recording immediately before placement by one of the  
113 researchers and stop recording immediately after collection. In both British and Chinese adults,  
114 the Axivity was set to capture triaxial acceleration data at 100 Hz with a dynamic range of  $\pm 8 g$ .

#### 115 *The Matrix 003*

116 The Matrix (██████████ XMatrix Tech. Co., Ltd, ██████████ China) was developed in 2021 to measure  
117 movement behaviours in the fifth wave of CHARLS (Zhao et al., 2014). It is a wrist-worn  
118 triaxial accelerometer similar to the Axivity, but it also has a gyroscope and a heart rate monitor.  
119 It is not possible to configure Matrix accelerometers to start and stop recording at specific times;  
120 its only options are to start immediately, after 24 hours, or after 48 hours. Therefore, in British  
121 adults, the Matrix was set to start recording just before 10am two working days after postal  
122 dispatch, whereas in Chinese adults it was set to start recording immediately before placement by  
123 one of the researchers. In both British and Chinese adults, the Matrix was set to capture triaxial  
124 acceleration data at 50 Hz with a dynamic range of  $\pm 8 g$ .

125 The Matrix was set to have a lower sampling rate than the Axivity (50 Hz versus 100 Hz) to  
126 ensure it had at least seven days of battery life whilst also capturing triaxial gyroscope data at 50  
127 Hz and heart rate data every 15 minutes. However, during data processing, both the Axivity and  
128 Matrix datasets were resampled using nearest neighbour interpolation at a rate of 50 Hz. Nearest  
129 neighbour interpolation resampling, recommended by Small et al. (2021), was used to avoid  
130 unintended smoothing of slower sampled data – specifically, the Matrix data – which has been  
131 shown to result in lower overall physical activity.

### 132 *Data processing*

133 Each participant's Matrix dataset was clipped to match their Axivity start and end times. Periods  
134 of Axivity non-wear time were then removed from the Matrix time series data and vice versa.  
135 Finally, gravity was calibrated in both accelerometers to ensure that at rest, the average  
136 magnitude of acceleration was 1 g (9.81 m/s<sup>2</sup>).

137 Participants were excluded if, for at least one of the accelerometers, the data could not be parsed,  
138 the device could not be calibrated, more than 1% of readings were 'clipped' (fell outside ±8 g for  
139 *biobankAccelerometerAnalysis*, ±3 g for *asleep*, and ±2 g for *stepcount*) before or after  
140 calibration, or the average acceleration was implausibly high (>100 mg for  
141 *biobankAccelerometerAnalysis* and *stepcount*, and >200 mg for *asleep*).

### 142 Primary outcomes (verification)

#### 143 Volume and intensity distribution of physical activity

144 The volume and intensity distribution of physical activity were derived from the Biobank  
145 Accelerometer Analysis Tool ([github.com/\[REDACTED\]biobankAccelerometerAnalysis](https://github.com/[REDACTED]/biobankAccelerometerAnalysis),  
146 v7.1.1), which was developed and validated by [REDACTED] (Walmsley et al.,

147 2021). Participants were excluded if they did not have sufficient wear time ( $\geq 3$  days in British  
148 and  $\geq 1$  day in Chinese adults, and data in each one-hour period of the 24-hour cycle), with non-  
149 wear time defined as unbroken episodes of at least 60 minutes during which the standard  
150 deviation (SD) of each axis of acceleration was less than 13 mg. To account for potential wear  
151 time diurnal bias, recording interruptions and non-wear time were imputed using the average  
152 values from the corresponding minute of the day on the remaining days of worn data.

153 The *biobankAccelerometerAnalysis* algorithm produced the following primary outcomes for each  
154 accelerometer separately: total wear time (days), average acceleration (mg per day), and time  
155 (hours per day) accumulated above incremental acceleration thresholds ( $>25$ - $200$  mg in 25-mg  
156 increments). Average acceleration refers to the Euclidean Norm Minus One (ENMO), a widely  
157 used measure of overall physical activity due to its correlation with physical activity energy  
158 expenditure (PAEE) (van Hees et al., 2013).

159 Secondary outcomes (analytical validation)

160 Time spent in 24-hour movement behaviours

161 The *biobankAccelerometerAnalysis* algorithm also produced the following secondary outcomes  
162 for each accelerometer separately: time (hours per day) spent in MVPA, in LPA, being  
163 sedentary, and in bed.

164 Sleep duration and efficiency

165 Sleep duration and efficiency were derived from a sleep staging algorithm

166 ([github.com/██████████/asleep](https://github.com/██████████/asleep), v0.4.13), which was also developed and validated by █████

167 █████ (Yuan et al., 2024). The *asleep* algorithm classifies each 30-second

168 epoch of acceleration data into one of the three sleep stages: 1) wake; 2) rapid eye movement

169 sleep (REM); and 3) non-rapid eye movement sleep (NREM). Participants were excluded if they  
170 did not have sufficient wear time ( $\geq 22$  hours per day for  $\geq 3$  days [including  $\geq 1$  weekend day] in  
171 British and  $\geq 1$  day in Chinese adults), with non-wear time defined as unbroken episodes of at  
172 least 90 minutes during which the SD of each axis of acceleration was less than 13 mg.

173 For each accelerometer separately, the *asleep* algorithm calculated the following sleep  
174 parameters for the longest sleep window over a noon-to-noon interval, with up to 60 minutes of  
175 sleep discontinuity allowed: overnight sleep duration (hours per night) and sleep efficiency  
176 (percentage of time in bed spent asleep).

177 Step count and cadence

178 Step count and cadence were derived from a hybrid machine learning and peak detection step  
179 counting algorithm ([github.com/stepcount](https://github.com/stepcount), v3.8.0), which was also developed and  
180 validated by [REDACTED] (Small et al., 2024). Participants were excluded if they  
181 did not have sufficient wear time ( $\geq 3$  days in British and  $\geq 1$  day in Chinese adults, and data in  
182 each one-hour period of the 24-hour cycle), with non-wear time defined as unbroken episodes of  
183 at least 90 minutes during which the SD of each axis of acceleration was less than 13 mg. To  
184 account for potential wear time diurnal bias, recording interruptions and non-wear time were  
185 imputed using the average values from the corresponding minute of the day on the remaining  
186 days of worn data.

187 The *stepcount* algorithm produced the following secondary outcomes for each accelerometer  
188 separately: overall daily step count (steps per day) and peak cadence (steps per minute). Overall  
189 daily step count was reported as the median number of steps taken per day across the monitoring

190 period. One-minute peak cadence was calculated as previously described by Saint-Maurice et al.  
191 (2020).

## 192 **Statistical analyses**

193 Descriptive statistics (mean [SD] where data were normally distributed, otherwise median [25<sup>th</sup>-  
194 75<sup>th</sup> percentile]) were calculated for all outcomes, with differences between British and Chinese  
195 adults being examined using the independent *t*-test or the Chi-squared test, and differences  
196 between the two accelerometer brands being examined using the paired *t*-test or the Wilcoxon  
197 signed-rank test.

198 We used 95% equivalence tests with a 10% equivalence zone to determine whether the 95%  
199 confidence interval (95% CI) for the mean of one accelerometer fell within  $\pm 10\%$  of the mean of  
200 the other accelerometer (Wellek, 2003). This was based on a previous study comparing three  
201 different accelerometers (Axivity AX3, ActiGraph GT9X, and GENEActiv) worn on both wrists  
202 (Rowlands et al., 2019), as well as many other validity studies of physical behaviour measures  
203 (O'Brien, 2021). Where data were not normally distributed, the log transformation of the original  
204 data were used for the equivalency analyses. As neither accelerometer can be considered the gold  
205 standard, equivalency analyses were carried out twice – i.e., with each accelerometer as the  
206 reference monitor. In all cases, equivalency was consistent regardless of the reference monitor.  
207 Therefore, results are presented with the Axivity as the reference monitor because it is more  
208 established.

209 The level of agreement between outputs from the two accelerometers was determined using  
210 intraclass correlation coefficients (ICCs; two-way mixed effects, absolute agreement, single  
211 measures) with 95% CI, and mean bias with 95% limits of agreement (LoA) (Bland & Altman,

212 1986). The ICC was classified as ‘poor’, ‘moderate’, ‘good’, or ‘excellent’ reliability if the lower  
213 band of the 95% CI of the ICC was <0.5, 0.5-0.75, >0.75-0.9, or >0.9, respectively (Koo & Li,  
214 2016).

215 Separate analyses were conducted for British and Chinese adults, and all analyses were  
216 conducted in RStudio (R version 4.4.1). Evidence of suggestive statistical significance was  
217 defined as  $p < 0.05$ .

## 218 **Results**

219 Of the 85 participants who were recruited in the UK, 82 had valid data for all outcomes from  
220 both accelerometers (52 [63.4%] female, mean [SD] age: 65.6 [10.2] years, 54 [65.9%] living in  
221 an urban area, mean [SD] BMI: 24.4 [3.6] kg/m<sup>2</sup>, 21 [25.6%] with mobility limitations). Of the  
222 117 participants who were recruited in China, 106 had valid data for all outcomes from both  
223 accelerometers. The Chinese participants had similar characteristics to the British participants,  
224 except they were more likely to have mobility limitations (64 [61.0%] female, mean [SD] age:  
225 66.7 [10.3] years, 58 [55.2%] living in an urban area, mean [SD] BMI: 24.4 [3.2] kg/m<sup>2</sup>, 49  
226 [46.7%] with mobility limitations). Descriptive statistics for all outcomes by nationality (British  
227 or Chinese) and accelerometer brand (Axivity or Matrix) are presented in Table 1.

228 In both British and Chinese adults, time spent in bed and sleeping were higher for the Matrix  
229 than the Axivity, whilst sleep efficiency and sedentary time were lower. In British adults, overall  
230 daily step count was lower for the Matrix than the Axivity, whereas the opposite was observed in  
231 Chinese adults. In British adults only, total wear time was higher for the Matrix than the Axivity,  
232 whilst overall physical activity (i.e., average acceleration), time accumulated above the lowest  
233 threshold of acceleration (25 mg), and time spent in LPA were lower. In Chinese adults only,

234 time accumulated above higher thresholds of acceleration ( $>100$  mg) were higher for the Matrix  
235 than the Axivity.

### 236 *Verification results*

237 Equivalency results for the volume and intensity distribution of physical activity by nationality  
238 are shown in Figure 1 and Supplementary Table S2. In Figure 1, markers are denoted in solid  
239 squares if the 95% CI for the mean of the Matrix 003 fell within  $\pm 10\%$  of the mean of the  
240 Axivity AX3, otherwise markers are denoted in hollow squares. In both British and Chinese  
241 adults, total wear time, average acceleration, and time accumulated above the two lowest  
242 thresholds of acceleration (25 and 50 mg) were equivalent between the two accelerometer brands  
243 (within  $\pm 10\%$  equivalence zone). Time accumulated above 75 mg was borderline equivalent in  
244 British ( $\pm 10.33\%$ ) but not equivalent in Chinese adults ( $\pm 14.33\%$ ), and time accumulated above  
245 the five highest thresholds of acceleration (100-200 mg) were not equivalent in either British or  
246 Chinese adults.

247 Agreement results for the volume and intensity distribution of physical activity by nationality are  
248 shown in Figure 2 and Supplementary Table S3. Bland-Altman plots are also shown in  
249 Supplementary Figure S2. In Figure 2, markers are denoted in solid squares if the lower band of  
250 the 95% CI of the ICC was good ( $>0.75$ ) or excellent ( $>0.9$ ), otherwise markers are denoted in  
251 hollow squares. In both British and Chinese adults, reliability between the two accelerometer  
252 brands was excellent for total wear time, average acceleration, and time accumulated above all  
253 thresholds of acceleration (25-200 mg).

254 *Analytical validation results*

255 Time spent in 24-hour movement behaviours

256 Equivalency results for time spent in 24-hour movement behaviours by nationality are shown in  
257 Figure 1 and Supplementary Table S2. In both British and Chinese adults, time spent in bed and  
258 being sedentary were equivalent between the two accelerometer brands (within  $\pm 10\%$   
259 equivalence zone); whereas time spent in MVPA was not equivalent, particularly in Chinese  
260 adults ( $\pm 17.74\%$  in British and  $\pm 28.20\%$  in Chinese adults). Time spent in LPA was equivalent  
261 between the two accelerometer brands in British ( $\pm 6.22\%$ ) but only borderline equivalent in  
262 Chinese adults ( $\pm 10.28\%$ ).

263 Agreement results for time spent in 24-hour movement behaviours are shown in Figure 2 and  
264 Supplementary Table S3. Bland-Altman plots are also shown in Supplementary Figure S2. In  
265 both British and Chinese adults, reliability between the two accelerometer brands was excellent  
266 for time spent in both LPA and MVPA but only moderate for time spent in bed. Reliability  
267 between the two accelerometer brands was excellent for sedentary time in British but only good  
268 in Chinese adults.

269 Sleep duration and efficiency

270 Equivalency results for sleep duration and efficiency are shown in Figure 1 and Supplementary  
271 Table S2. In both British and Chinese adults, both overnight sleep duration and sleep efficiency  
272 were equivalent between the two accelerometer brands (within  $\pm 10\%$  equivalence zone).

273 Agreement results for sleep duration and efficiency are shown in Figure 2 and Supplementary  
274 Table S3. Bland-Altman plots are also shown in Supplementary Figure S2. In both British and  
275 Chinese adults, reliability between the two accelerometer brands was good for sleep efficiency.

276 Reliability between the two accelerometer brands was good for overnight sleep duration in  
277 British but poor in Chinese adults.

278 Step count and cadence

279 Equivalency results for step count and cadence are shown in Figure 1 and Supplementary Table  
280 S2. In both British and Chinese adults, peak cadence was equivalent between the two  
281 accelerometer brands (within  $\pm 10\%$  equivalence zone). Overall daily step count was also  
282 equivalent between the two accelerometer brands in Chinese ( $\pm 2.93\%$ ) but only borderline  
283 equivalent in British adults ( $\pm 10.50\%$ ).

284 Agreement results for step count and cadence are shown in Figure 2 and Supplementary Table  
285 S3. Bland-Altman plots are also shown in Supplementary Figure S2. In both British and Chinese  
286 adults, reliability between the two accelerometer brands was excellent for both overall daily step  
287 count and peak cadence.

## 288 **Discussion**

289 In both British and Chinese adults, average acceleration was equivalent between the Axivity  
290 AX3 and the Matrix 003. The vast majority of machine learning-derived physical activity and  
291 sleep outcomes were also equivalent or borderline equivalent. The only exceptions were time  
292 spent above higher thresholds of acceleration ( $>75\text{ mg}$ ), including time spent in MVPA.

293 However, there were a few notable cross-nationality differences. In Chinese compared to British  
294 adults, equivalency between the two accelerometer brands was lower for time spent in both LPA  
295 and MVPA but higher for overall daily step count. Furthermore, the ICCs were good or excellent  
296 ( $>0.75$ ) for all outcomes, except for time spent in bed in both British and Chinese adults and  
297 overnight sleep duration in Chinese adults only.

298 We found that average acceleration was equivalent between the Axivity AX3 and the Matrix 003  
299 when measured at the dominant wrist. In contrast, in 56 young British adults (mean age 24.5  
300 years), Rowlands et al. (2019) found that average acceleration measured at the dominant wrist  
301 was approximately 10% higher for the Axivity AX3 and the GENEActiv than for the ActiGraph  
302 GT9X. Our results suggest that accelerometry data from the UK Biobank, China Kadoorie  
303 Biobank, ELSA, and CHARLS can be pooled and/or compared because they all used the Axivity  
304 AX3 or the Matrix 003 on the dominant wrist. We also found that time spent above lower  
305 thresholds of acceleration ( $\leq 75$  mg) were equivalent between the Axivity AX3 and the Matrix  
306 003, but time spent above higher thresholds ( $> 75$  mg) were not. In contrast, Rowlands et al.  
307 (2019) found that the intensity gradient – a single metric describing the distribution of  
308 acceleration intensity across the 24-hour day (Rowlands et al., 2018) – was equivalent  
309 irrespective of accelerometer brand or wrist. The lack of equivalency in our study was most  
310 likely due to the 10% equivalence zone being too strict for less common activities that tend to  
311 have low magnitudes and/or high variability (Rowlands et al., 2019). This issue may have been  
312 exacerbated by the particularly low levels of higher-intensity physical activity typically observed  
313 in older adults.

314 Our findings, together with those of Rowlands et al. (2019), have implications for harmonising  
315 data across large-scale epidemiological cohorts that have used wrist-worn accelerometers.  
316 Historically, however, many studies have relied on waist-worn devices. In a systematic review of  
317 observational studies measuring device-based physical behaviours in adults, Pulsford et al.  
318 (2023) reported that the waist was the most common wear location, used in 53% of study waves  
319 compared to 20% for the wrist. Nonetheless, wear time compliance was higher for wrist-worn  
320 devices than for waist-worn ones, supporting the increasing preference for wrist placement in

321 more recent studies. Furthermore, some studies used thigh-worn accelerometers (5% of study  
322 waves), which enable more accurate determination of posture and stepping, while others used  
323 multiple wear locations (13% of study waves) to measure different dimensions of physical  
324 behaviour, such as intensity, posture, and biological state (e.g., asleep or awake).

325 Substantial differences in accelerometer outputs have been observed between wear locations  
326 under free-living conditions. For example, a meta-analysis by Gall, Sun, and Smuck (2022)  
327 reported that wrist-worn accelerometers recorded, on average, 3,537 more steps per day than  
328 waist-worn ones. Wrist-worn devices also captured more time in MVPA and less sedentary time  
329 compared to waist-worn ones. Similarly, Maylor et al. (2023) found that the wrist-worn Axivity  
330 AX3 recorded higher average acceleration and a wider range of acceleration values than the  
331 thigh-worn activPAL micro.

332 Therefore, we are currently conducting two new free-living validation studies – in adults aged  
333 18-30 years and those aged 40 years or older – to establish whether the same machine learning  
334 algorithms can be applied to different accelerometer wear locations, or if adjustment factors or  
335 location-specific algorithms are needed for data harmonisation.

336 To our knowledge, this is the first study to show that the same machine learning algorithms can  
337 be applied to different brands of accelerometer in different populations. However, the cross-  
338 nationality differences observed support the findings of two previous free-living validation  
339 studies in British and Chinese adults: CAPTURE-24 (Walmsley et al., 2021) and CAPTURE-  
340 24CN (Chen et al., 2023), respectively. In CAPTURE-24, 82% of all MVPA instances involved  
341 walking or cycling because almost all of the participants were office workers. In contrast, in  
342 CAPTURE-24CN, only 42% of MVPA instances involved walking or cycling, while the rest  
343 consisted of farm or construction work (Chen et al., 2023). Therefore, our finding that

344 equivalency between the two accelerometer brands was lower for time spent in both LPA and  
345 MVPA but higher for overall daily step count in Chinese compared to British adults may be due  
346 to the UK-derived machine learning algorithms being better at identifying walking and cycling  
347 than farm or construction work.

348 Reprocessing the raw acceleration data from Chinese adults using a China-derived machine  
349 learning algorithm (Chen et al., 2023) improved equivalency between the two accelerometer  
350 brands for time spent in MVPA, bringing it closer to the level of equivalency observed in British  
351 adults (Supplementary Table S4). However, it still fell outside the predefined 10% equivalence  
352 zone, and equivalency for time spent in LPA was not improved by the China-derived algorithm.  
353 These findings highlight the need for further algorithm retraining using the additional wearable  
354 camera data collected in the current study.

355 It was somewhat surprising that equivalency between the two accelerometer brands for overall  
356 daily step count was higher in Chinese than in British adults, given that the *stepcount* machine  
357 learning algorithm was trained on data from British adults. Koffman and Muschelli (2024)  
358 applied five open-source and one proprietary algorithm to three publicly available datasets with  
359 ground truth step counts, including the free-living OxWalk dataset used to train the *stepcount*  
360 algorithm. They found that the *stepcount* algorithm was the most accurate, achieving the highest  
361 F1 score ( $0.89 \pm 0.11$ ) and the lowest mean absolute percentage error ( $8.6 \pm 9.0\%$ ). However, the  
362 adults in OxWalk were much younger than those in the current study (mean age: 38.5 years vs.  
363 65.6 and 66.7 years), highlighting the need for additional free-living validation studies across a  
364 wider range of ages and nationalities.

365 The sleep outcomes had the lowest ICCs in the current study, with the Matrix 003 recording  
366 more time in bed, higher overnight sleep duration, and poorer sleep efficiency than the Axivity

367 AX3. This was particularly the case in Chinese adults, which may have been due to the shorter  
368 monitoring period of only two nights. Yuan et al. (2024) previously found that at least three days  
369 were required for stable weekly sleep parameter estimates. Further investigation using the  
370 additional sleep diary data collected is needed before pooling or comparing sleep outcomes from  
371 these two different accelerometer brands.

### 372 *Strengths and limitations*

373 A strength of this study was the concurrent wear of two different brands of accelerometer (the  
374 well-established Axivity AX3 and the newly developed Matrix 003) in a random positional order  
375 on the dominant wrist 24 hours per day for up to eight days. Furthermore, the raw acceleration  
376 data from both accelerometer brands were processed identically using three well-validated, open-  
377 source machine learning algorithms to extract meaningful physical activity and sleep outcomes  
378 (Walmsley et al., 2021; Small et al., 2024; Yuan et al., 2024). Another strength of this study was  
379 that it was conducted in both the UK and China, allowing us to identify some important cross-  
380 nationality differences. The same methods were used in both countries, except the  
381 accelerometers were distributed and returned via post in the UK but in person in China. The  
382 monitoring period was also shorter in China than in the UK (two nights and one day versus eight  
383 days), meaning the findings may be less robust. Another limitation of this study was the use of  
384 self-selected cohorts, which may limit the generalisability of the findings. Moreover, older adults  
385 may exhibit substantially different physical activity and sleep patterns compared to middle-aged  
386 and younger adults.

## 387 *Conclusions*

388 In both British and Chinese adults, average acceleration – a widely used measure of overall  
389 physical activity – was equivalent between the Axivity AX3 and the Matrix 003 worn in a  
390 random positional order on the dominant wrist. This finding suggests that data from these two  
391 different brands of accelerometer can be pooled and/or compared, including across large-scale  
392 epidemiological cohorts such as the UK Biobank, China Kadoorie Biobank, ELSA, and  
393 CHARLS. The vast majority of machine learning-derived physical activity and sleep outcomes  
394 were also equivalent or borderline equivalent. However, the cross-nationality differences  
395 observed likely stem from the UK-developed algorithms being better at identifying walking and  
396 cycling than farm or construction work. This highlights the need for further population-specific  
397 algorithm development and validation.

## 398 **Funding**

399 LB was supported by both the Economic and Social Research Council (ES/T014091/1) and the  
400 Wellcome Trust (223100/Z/21/Z), and MB was supported by the Economic and Social Research  
401 Council (ES/T014091/1). AD's research team is supported by a range of grants from the  
402 Wellcome Trust (223100/Z/21/Z, 227093/Z/23/Z), Novo Nordisk, Swiss Re, Boehringer  
403 Ingelheim, National Institutes of Health's Oxford-Cambridge Scholars Program, EPSRC Centre  
404 for Doctoral Training in Health Data Science (EP/S02428X/1), British Heart Foundation Centre  
405 of Research Excellence (RE/18/3/34214), and funding administered by the Danish National  
406 Research Foundation in support of the Pioneer Centre for SMARTbiomed. The funders had no  
407 role in conceptualisation, design, data collection, analysis, decision to publish, or preparation of  
408 the manuscript. For the purpose of open access, the authors have applied a Creative Commons  
409 Attribution (CC BY) licence to any Authors Accepted Manuscript version arising.

410 **Data availability**

411 The dataset analysed during the current manuscript is not publicly available due to institutional  
412 data sharing restrictions, but we anticipate anonymised data to be shared in future as part of a  
413 wider project.

414 **References**

415 Bland, J. M., & Altman, D. G. (1986). Statistical methods for assessing agreement between two  
416 methods of clinical measurement. *The Lancet*, 327(8476), 307–310.

417 [https://doi.org/10.1016/S0140-6736\(86\)90837-8](https://doi.org/10.1016/S0140-6736(86)90837-8)

418 Bull, F. C., Al-Ansari, S. S., Biddle, S., Borodulin, K., Buman, M. P., Cardon, G., ... Willumsen,  
419 J. F. (2020). World Health Organization 2020 guidelines on physical activity and sedentary  
420 behaviour. *British Journal of Sports Medicine*, 54(24), 1451–1462.

421 <https://doi.org/10.1136/bjsports-2020-102955>

422 Chaput, J. P., Dutil, C., Featherstone, R., Ross, R., Giangregorio, L., Saunders, T. J., ... Carrier,  
423 J. (2020). Sleep duration and health in adults: an overview of systematic reviews. *Applied*  
424 *Physiology, Nutrition, and Metabolism*, 45(10 Suppl. 2), S218–S231.

425 <https://doi.org/10.1139/apnm-2020-0034>

426 Chen, Y., Chan, S., Bennett, D., Chen, X., Wu, X., Ke, Y., ... Doherty, A. (2023). Device-  
427 measured movement behaviours in over 20,000 China Kadoorie Biobank participants.  
428 *International Journal of Behavioral Nutrition and Physical Activity*, 20, 138.

429 <https://doi.org/10.1186/s12966-023-01537-8>

430 Doherty, A., Jackson, D., Hammerla, N., Plötz, T., Olivier, P., Granat, M. H., ... Wareham, N. J.  
431 (2017). Large scale population assessment of physical activity using wrist worn accelerometers:  
432 The UK Biobank Study. *PLoS ONE*, *12*(2), e0169649.  
433 <https://doi.org/10.1371/journal.pone.0169649>

434 Ekelund, U., Tarp, J., Steene-Johannessen, J., Hansen, B. H., Jefferis, B., Fagerland, M. W., ...  
435 Lee, I. M. (2019). Dose-response associations between accelerometry measured physical activity  
436 and sedentary time and all cause mortality: systematic review and harmonised meta-analysis.  
437 *BMJ*, *366*, 14570. <https://doi.org/10.1136/bmj.14570>

438 Gall, N., Sun, R., & Smuck, M. (2022). A comparison of wrist- versus hip-worn ActiGraph  
439 sensors for assessing physical activity in adults: A systematic review. *Journal for the*  
440 *Measurement of Physical Behaviour*, *5*(4), 252-262. <https://doi.org/10.1123/jmpb.2021-0045>

441 Goldsack, J. C., Coravos, A., Bakker, J. P., Bent, B., Dowling, A. V., Fitzer-Attas, C., ... Dunn,  
442 J. (2020). Verification, analytical validation, and clinical validation (V3): The foundation of  
443 determining fit-for-purpose for Biometric Monitoring Technologies (BioMeTs). *NPJ Digital*  
444 *Medicine*, *3*, 55. <https://doi.org/10.1038/s41746-020-0260-4>

445 Koo, T. K., & Li, M. Y. (2016). A guideline of selecting and reporting intraclass correlation  
446 coefficients for reliability research. *Journal of Chiropractic Medicine*, *15*(2), 155–163.  
447 <https://doi.org/10.1016/j.jcm.2016.02.012>

448 Maylor, B. D., Edwardson, C. L., Clarke-Cornwell, A. M., Davies, M. J., Dawkins, N. P.,  
449 Dunstan, D. W., ... Rowlands, A. V. (2023). Physical activity assessed by wrist and thigh worn  
450 accelerometry and associations with cardiometabolic health. *Sensors*, *23*(17), 7353.  
451 <https://doi.org/10.3390/s23177353>

452 Menai, M., van Hees, V. T., Elbaz, A., Kivimäki, M., Singh-Manoux, A., & Sabia, S. (2017).  
453 Accelerometer assessed moderate-to-vigorous physical activity and successful ageing: Results  
454 from the Whitehall II study. *Scientific Reports*, 8, 45772. <https://doi.org/10.1038/srep45772>

455 O'Brien, M. W. (2021). Implications and recommendations for equivalence testing in measures  
456 of movement behaviours: a scoping review. *Journal for the Measurement of Physical Behaviour*,  
457 4(4), 353-362. <https://doi.org/10.1123/jmpb.2021-0021>

458 Pulsford, R. M., Brocklebank, L., Fenton, S. A. M., Bakker, E., Mielke, G. I., Tsai, L., ...  
459 Stamatakis, E. (2023). The impact of selected methodological factors on data collection  
460 outcomes in observational studies of device-measured physical behaviour in adults: A systematic  
461 review. *International Journal of Behavioural Nutrition and Physical Activity*, 20(26).  
462 <https://doi.org/10.1186/s12966-022-01388-9>

463 Rowlands, A. V. (2018). Moving forward with accelerometer-assessed physical activity: Two  
464 strategies to ensure meaningful, interpretable, and comparable measures. *Pediatric Exercise*  
465 *Science*, 30(4), 450–456. <https://doi.org/10.1123/pes.2018-0201>

466 Rowlands, A. V., Plekhanova, T., Yates, T., Mirkes, E. M., Davies, M., Khunti, K., &  
467 Edwardson, C. L. (2019). Providing a basis for harmonization of accelerometer-assessed  
468 physical activity outcomes across epidemiological datasets. *Journal for the Measurement of*  
469 *Physical Behaviour*, 2(3), 131–142. <https://doi.org/10.1123/jmpb.2018-0073>

470 Saint-Maurice, P. F., Troiano, R. P., Bassett Jr, D. R., Graubard, B. I., Carlson, S. A., Shiroma,  
471 E. J., ... Matthews, C. E. (2020). Association of daily step count and step intensity with mortality  
472 among US adults. *JAMA*, 323(12), 1151–1160. <https://doi.org/10.1001/jama.2020.1382>

473 Small, S., Khalid, S., Dhiman, P., Chan, S., Jackson, D., Doherty, A., & Price, A. (2021). Impact  
474 of Reduced Sampling Rate on Accelerometer-Based Physical Activity Monitoring and Machine  
475 Learning Activity Classification. *Journal for the Measurement of Physical Behaviour*, 4(4), 298–  
476 310. <https://doi.org/10.1123/jmpb.2020-0061>

477 Small, S. R., Chan, S., Walmsley, R., von Fritsch, L., Acquah, A., Mertes, G., ... Doherty, A.  
478 (2024). Self-supervised machine learning to characterize step counts from wrist-worn  
479 accelerometers in the UK Biobank. *Med Sci Sports Exerc*, 56(10), 1945-1953.  
480 <https://doi.org/10.1249/MSS.0000000000003478>

481 Steptoe, A., Breeze, E., Banks, J., & Nazroo, J. (2013). Cohort profile: the English longitudinal  
482 study of ageing. *International Journal of Epidemiology*, 42(6), 1640–1648.  
483 <https://doi.org/10.1093/ije/dys168>

484 van Hees, V. T., Gorzelniak, L., Dean León, E. C., Eder, M., Pias, M., Taherian, S., ... Brage, S.  
485 (2013). Separating movement and gravity components in an acceleration signal and implications  
486 for the assessment of human daily physical activity. *PLoS One*, 8(4), e61691.  
487 <https://doi.org/10.1371/journal.pone.0061691>

488 van Hees, V. T., Fang, Z., Langford, J., Assah, F., Mohammad, A., da Silva, I. C. M., ... Brage,  
489 S. (2014). Autocalibration of accelerometer data for free-living physical activity assessment  
490 using local gravity and temperature: An evaluation on four continents. *J Appl Physiol (1985)*,  
491 117(7), 738–744. <https://doi.org/10.1152/jappphysiol.00421.2014>

492 Walmsley, R., Chan, S., Smith-Byrne, K., Ramakrishnan, R., Woodward, M., Rahimi, K., ...  
493 Doherty, A. (2021). Reallocation of time between device-measured movement behaviours and

494 risk of incident cardiovascular disease. *British Journal of Sports Medicine*, 56(18), 1008-1017.  
495 <https://doi.org/10.1136/bjsports-2021-104050>

496 Wasfy, M. M., & Lee, I. M. (2022). Examining the dose-response relationship between physical  
497 activity and health outcomes. *NEJM Evidence*, 1(12), EVIDra2200190.  
498 <https://doi.org/10.1056/EVIDra2200190>

499 Wellek, S. (2003). *Testing statistical hypotheses of equivalence*. Chapman & Hall/CRC

500 Yuan, H., Plekhanova, T., Walmsley, R., Reynolds, A. C., Maddison, K. J., Bucan, M., ...  
501 Doherty, A. (2024). Self-supervised learning of accelerometer data provides new insights for  
502 sleep and its association with mortality. *NPJ Digital Medicine*, 7(1), 86.  
503 <https://doi.org/10.1038/s41746-024-01065-0>

504 Zhao, Y., Hu, Y., Smith, J. P., Strauss, J., & Yang, G. (2014). Cohort profile: the China Health  
505 and Retirement Longitudinal Study (CHARLS). *International Journal of Epidemiology*, 43(1),  
506 61–68. <https://doi.org/10.1093/ije/dys203>

Table 1. Physical activity and sleep outcomes by nationality and accelerometer brand.

	British adults (N=82)		<i>P</i> -value	Chinese adults (N=106)	
	Axivity AX3	Matrix 003		Axivity AX3	Matrix 003
Total wear time (days)	6.8 (6.6-7.0)	6.9 (6.6-7.0)	<b>0.041</b>	1.6 (1.5-1.6)	1.6 (1.5-1.6)
Average acceleration (mg)	27.9 (23.5-34.5)	26.8 (21.7-32.4)	<b>&lt;0.001</b>	25.8 (19.7-34.6)	26.1 (19.2-32.7)
Time above 25 mg (hours/day)	6.9 (1.7)	6.8 (1.7)	<b>0.001</b>	7.5 (6.0-9.2)	7.7 (5.8-8.9)
Time above 50 mg	4.3 (1.4)	4.2 (1.4)	0.162	4.1 (3.0-5.5)	4.4 (2.8-5.7)
Time above 75 mg	2.6 (1.1)	2.6 (1.1)	0.277	2.3 (1.2-3.4)	2.3 (1.3-3.4)
Time above 100 mg	1.5 (1.1-2.2)	1.6 (1.0-2.3)	0.594	1.2 (0.6-2.0)	1.3 (0.6-2.0)
Time above 125 mg	1.0 (0.6-1.5)	1.0 (0.6-1.4)	0.776	0.7 (0.3-1.2)	0.6 (0.3-1.3)
Time above 150 mg	0.6 (0.4-1.0)	0.6 (0.3-1.0)	0.876	0.3 (0.1-0.8)	0.4 (0.2-0.8)
Time above 175 mg	0.4 (0.2-0.7)	0.4 (0.2-0.7)	0.746	0.2 (0.1-0.5)	0.2 (0.1-0.5)
Time above 200 mg	0.3 (0.1-0.5)	0.2 (0.1-0.5)	0.495	0.1 (0.0-0.4)	0.2 (0.0-0.4)
Time in bed (hours/day)	7.7 (7.1-8.3)	7.8 (7.3-8.5)	<b>&lt;0.001</b>	7.0 (6.3-8.2)	7.3 (6.6-8.7)
SB	10.8 (1.8)	10.6 (1.7)	<b>&lt;0.001</b>	10.6 (2.8)	10.3 (2.5)
LPA	4.4 (3.5-5.0)	4.3 (3.4-5.1)	<b>0.001</b>	4.5 (3.4-6.2)	4.6 (3.2-6.2)
MVPA	1.1 (0.6-1.6)	1.0 (0.6-1.6)	0.077	1.0 (0.5-1.9)	1.0 (0.4-1.9)
Sleep duration (hours/night)	5.3 (3.8-6.3)	6.0 (3.9-6.9)	<b>0.001</b>	5.5 (4.3-6.3)	6.1 (5.3-6.8)
Sleep efficiency (%)	75.1 (58.5-83.8)	70.7 (48.4-79.9)	<b>&lt;0.001</b>	85.9 (78.1-92.3)	83.3 (73.9-88.3)
Step count (steps/day)	11,018 (4,138)	10,726 (3,995)	<b>&lt;0.001</b>	10,409 (7,686-15,903)	10,534 (7,005-15,343)
Peak cadence (steps/minute)	116 (13)	116 (14)	1.000	107 (14)	106 (14)

Mean (SD) is presented where data were normally distributed, otherwise median (25<sup>th</sup>-75<sup>th</sup> percentile).

mg = milli-gravity units (dynamic acceleration expressed as ENMO).

SB, sedentary behaviour; LPA, light physical activity; MVPA, moderate-vigorous physical activity; SD, standard deviation; ENMO, Euclidean Norm Minus One.

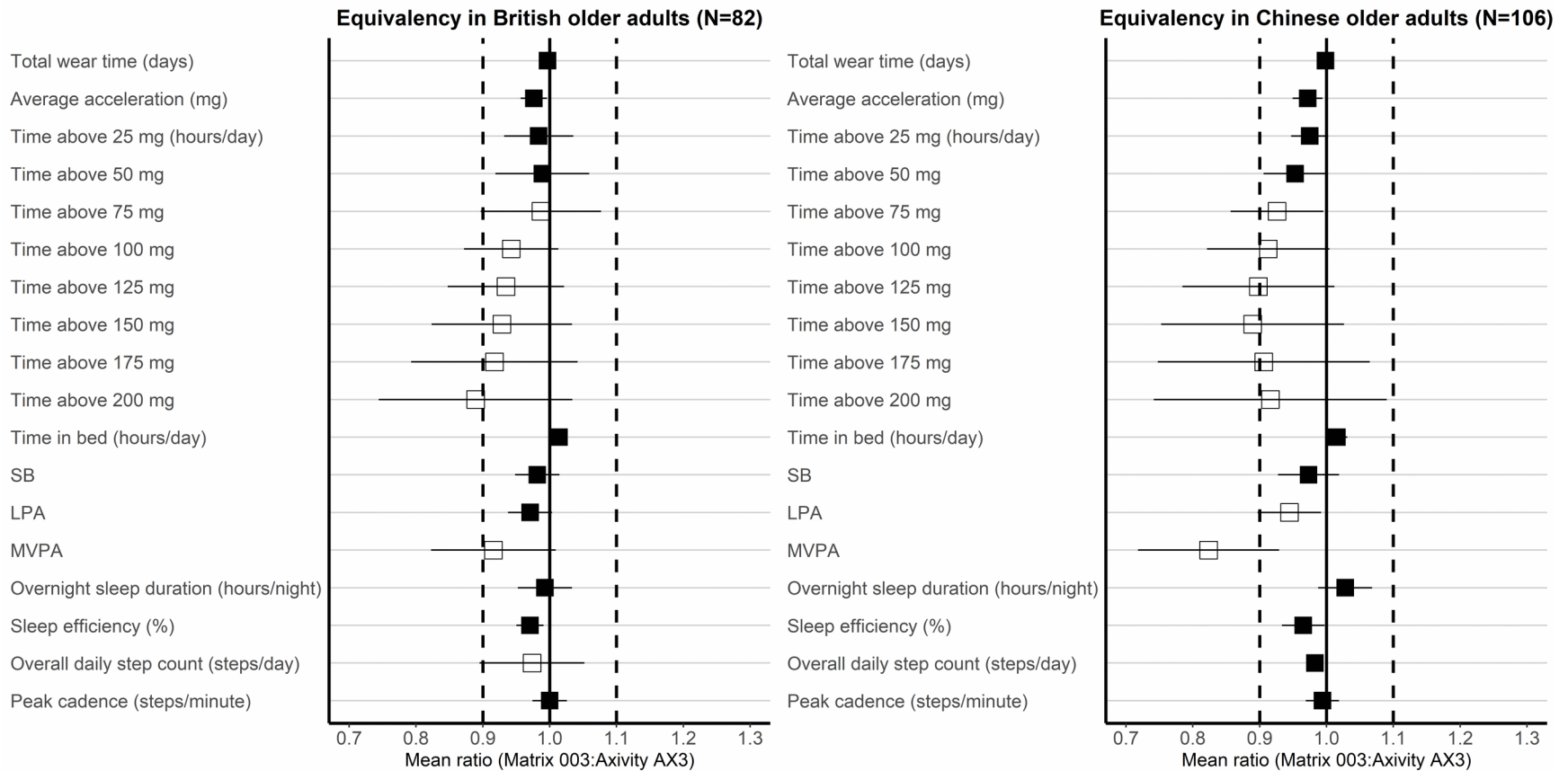


Figure 1. Equivalency between accelerometer brands for the physical activity and sleep outcomes by nationality.

For equivalency, the dashed lines indicate the a priori 10% equivalence zone. The horizontal black lines show 95% CIs. Markers are denoted in solid squares if the 95% CI for the mean of the Matrix 003 accelerometer fell within  $\pm 10\%$  of the mean of the Axivity AX3 accelerometer, otherwise markers are denoted in hollow squares.

mg = milli-gravity units (dynamic acceleration expressed as ENMO).

SB, sedentary behaviour; LPA, light physical activity; MVPA, moderate-vigorous physical activity; 95% CI, 95% confidence interval; ENMO, Euclidean Norm Minus One.

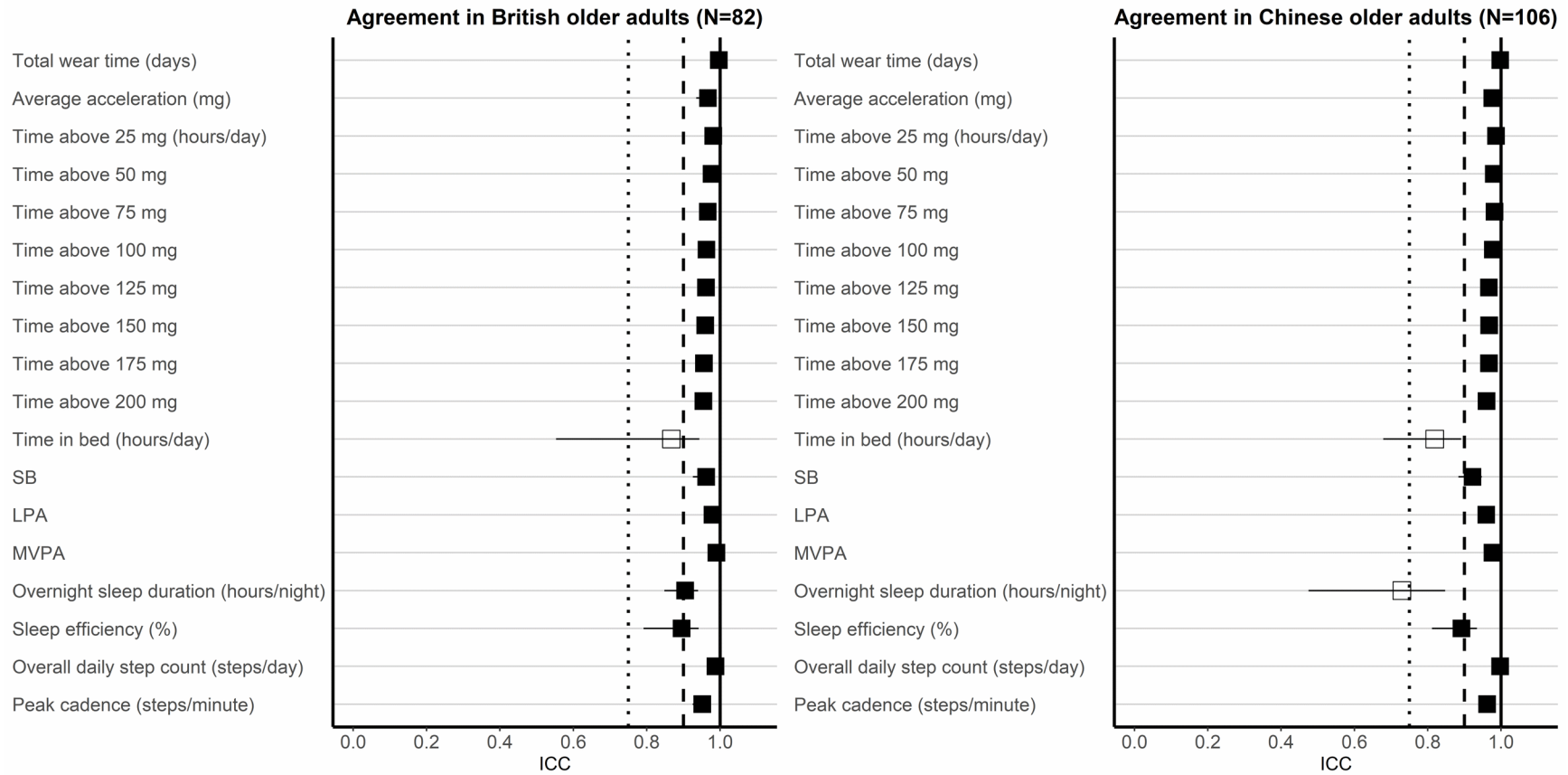


Figure 2. Agreement between accelerometer brands for the physical activity and sleep outcomes by nationality.

For agreement, the dotted line represents ICCs of 0.75 (good reliability) and the dashed line represents ICCs of 0.9 (excellent reliability). The horizontal black lines show 95% CIs. Markers are denoted in solid squares if the lower band of the 95% CI of the ICC was >0.75, otherwise markers are denoted in hollow squares.

mg = milli-gravity units (dynamic acceleration expressed as ENMO).

SB, sedentary behaviour; LPA, light physical activity; MVPA, moderate-vigorous physical activity; ICC, intraclass correlation coefficient; 95% CI, 95% confidence interval; ENMO, Euclidean Norm Minus One.