1

2

# Validity of Optical Heart Rate Measurement in Commercially Available

# Wearable Fitness Tracking Devices

5

6

7   Jason Thomas[1]*, Patrick Doyle[1☐], J. Andrew Doyle[1]

8

9

10

11   [1] Department of Kinesiology and Health, College of Education and Human Development, Georgia State

12   University, Atlanta, Georgia, United States of America

13   [☐] Current Address: Institute for Disaster Management, College of Public Health, University of Georgia,

14   Athens, Georgia, United States of America

15

16   *  Corresponding author

17   E-mail: jthomas31@gsu.edu

18

## Abstract

**Background:** Wearable fitness tracking devices have risen in popularity for athletes and the general population and are increasingly integrated into smartwatch technology. Many devices incorporate optical heart rate (HR) measurement by photoplethysmography which provides data used to monitor and track exercise training intensities, progress, and other health and fitness related parameters.

**Objective:** To determine the validity of optical heart rate measurement in three fitness tracking devices while resting, walking, and running.

**Methods:** Twenty subjects (10 male, 10 female) completed the research study based on the ANSI/CTI standards for physical activity monitoring of heart rate under 4 different conditions: sedentary (SED), treadmill walking (WLK), running (RUN), and dynamic running/walking (DYN). Subjects wore 3 optical heart rate devices: Polar OH1 (OH1) on the right forearm, Apple Watch 4 (AW4) on the right wrist and Garmin Forerunner 945 (FR945) on the left wrist. A Polar H10 (H10), a chest strap device, was the criterion HR measurement device. SED, WLK, and RUN were all 7-minute protocols with 1 minute of standing, 5 minutes of prescribed activity, and 1 final minute of standing. The DYN protocol was a 12-minute protocol with 1 minute of standing, 10 minutes of variable intensity walking and running, and 1 minute of standing. Raw HR data was extracted from each device and temporally aligned with the criterion H10 HR data for analysis.

**Results:**

The mean absolute deviation (MAD, measured in beats per minute) for the three experimental devices (OH1, AW4, FR945, respectively) for SED was 1.31, 1.33, and 2.03; for WLK was 2.79, 2.58, and 5.19; for RUN were 4.00, 4.29, and 6.51; and for DYN was 2.60, 2.44, and 2.44. The mean absolute percent error (MAPE) for the three experimental devices (OH1, AW4, FR945, respectively) for SED was 1.78%, 1.89%,

2

41    and 2.81%; for WLK was 3.15%, 3.18%, and 5.93%; for RUN was 3.43%, 3.51%, and 5.25%; and for DYN

42    was 2.05%, 1.95%, and 5.47%.  The intraclass correlation for each device across all conditions was .991

43    (OH1), .984 (AW4), and .697 (FR945).

44    **Conclusions:** At rest, and during both steady-state and variable-speed treadmill walking and running, the

45    Polar OH1, Garmin Forerunner 945, and Apple Watch 4 optical HR monitors demonstrated a level of

46    accuracy well within that required by the ANSI/CTA Standard (2018) for physical activity monitoring

47    devices for heart rate measurement (i.e., <10% Mean Absolute Percent Error).  Therefore, consumers

48    can have confidence that these devices provide HR data with accuracy that conforms to the

49    performance criteria recommended for consumer electronics.

50    **Keywords:** photoplethysmography; heart rate monitor; smartwatch; fitness tracking device

51    ## Introduction

52    ## Background

53    Wearable fitness tracking devices have risen in popularity over the past decade and have been the top

54    fitness trend numerous years while approaching nearly a $100 billion industry (1).  These devices were

55    initially developed as either rudimentary mechanical pedometers attached to a shoe or waist band, or

56    electrode chest strap heart rate monitors that are often deemed uncomfortable and cumbersome.  As

57    technology has advanced, wearable fitness devices have integrated improved technologies including

58    GPS, accelerometers, altimeters, and photosensors.  Further, they are increasingly integrated into more

59    user-friendly and comfortable devices, specifically wrist-worn watches, and arm bands.  As heart rate

60    (HR) monitoring is arguably the key component of fitness monitoring, a principle technological advance

61    has been the integration of photoplethysmography (PPG), which uses a light emitting diode and

62    photosensor to measure microvascular blood volume changes which is consequently associated with

63    heart rate (2).

64    The advancement and integration of PPG technology into wrist-worn devices has granted the end-user

65    with a wealth of information including caloric expenditure, oxygen consumption ($VO_2$), heart rate

66    variability, sleep patterns, recovery, and training intensity.  All this information provided to users is

67    based on manufacturer-specific algorithms computed from heart rate collected via PPG technology.

68    Therefore, the validity of the heart rate measurement from these PPG devices is of key importance.

69    Several studies have been completed to assess the validity of a variety of activity tracking devices which

70    use the PPG technology.  Although several devices, including the OH1, Apple Watch series, and Garmin

71    Forerunner series, have been deemed valid, the results of the studies must be interpreted narrowly as

72    various methodological differences or concerns exist between studies.  As device availability has grown

73    immensely and rapidly, the current body of research lacks results that can more confidently discern the

74    validity of the devices across the general population.

75    The OH1 has been previously studied and was deemed valid for moderate and high intensity activities

76    (3, 4),  appeared more valid compared to a wrist-based device by the same manufacturer(5), and

77    showed decreased validity with arm-based activities (i.e. tennis) (6).  A key limitation of these studies is

78    the application of the study results to the wider population as the studies lacked balanced diversity in

79    either BMI, skin tone, or sex.  Likewise, studies using the Apple Watch series have suggested device

80    validity, but different methodological issues exist.  The methodological concerns were comprised of

81    various issues such as  recording heart rate from a single timepoint (7), using a model of tachycardia (8),

82    implementing a single subject design (9), or failing to report key validity metrics such as MAD, MAPE,

83    and ICC (10).  Additionally, these studies also lacked the diversity in key subject demographics, similar to

84    the limitations with OH1 research.  The Garmin Forerunner series has a very limited amount of

85    information available in the literature.  The existing data has either suggested poor validity in prior

86    versions to the FR945 (7) or has suffered from methodological issues related to heart rate recording

87    frequency (11).

4

88    Aside from specific device validity, the current body of research for all PPG activity tracking devices

89    suffers from numerous methodological differences that limits our ability to apply the results to the

90    general population.  Existing studies generally lack cohesion between different exercise types,

91    intensities, and durations.  Some studies have been completed to assess the validity of a single device

92    across multiple exercise modalities (4), while others have researched numerous devices across a variety

93    of intensities and exercise modalities, but with shorter data recording times (12).   More recent studies

94    have investigated multiple devices, but intensity was extremely high and duration extremely short (13).

95    Other studies have utilized different modes of exercise but lacked varying levels of intensity within the

96    modes (14, 15).  A recent study implemented activity modes and intensity with a better variation, but

97    only incorporated a single demographic (Caucasian) in the subject group (10)

98    The lack of variation in subject demographics is visible across many studies.  Variations in skin tone, BMI,

99    sex, and age have been suggested as potential confounding factors to proper validity testing for PPG

100    technology.  Variations in skin tone appears to affect validity as use of a typical green-light LED diode

101    (often integrated into many devices) has resulted in a 1.04 BPM error rate in light-skinned individuals,

102    and as much as a 10.9 BPM error in dark-skinned individuals (16).  There is also evidence to suggest that

103    as BMI increases, PPG waveform can change as much as 43% between obese and non-obese individuals

104    (17).  Some studies have attempted to address these concerns but have had limitations.  For instance, a

105    recent study did investigate skin tone and PPG using an Apple Watch, but the study subjects only

106    represented 3 of the 6 Fitzpatrick skin tone designations (18).  Another study had all 6 skin tones

107    represented, but only had 10 subjects total such that certain skin tones were only represented by a

108    single subject (19).  Additionally, few studies have specifically recruited subjects to represent a diversity

109    in BMI or gender.

110    Recently, the American National Standards Institute (ANSI) and Consumer Technology Association (CTA)

111    developed the ANSI/CTA standards for investigating the validity and reliability of consumer electronic

5

112   fitness devices.  These standards provide a consistent, balanced, and equitable basis for subject

113   selection and activity parameters so that consumer devices can be evaluated in a standardized manner.

114   The activity parameters outline optimal intensity levels and duration for different modes of activity.

115   Additionally, subject selection requirements ensure a diverse population relative to age, gender, body

116   mass, and skin tone or complexion.

117   Study Objective

118   PPG technology is being widely implemented to determine HR in an increasing number of devices to

119   appeal to a broader market of consumers globally.  As such, it is important to determine if the existing

120   device validity evidenced by previous studies is representative of a diverse population and activities or if

121   the results can only be applied to the limited subject demographics and activities of the respective

122   studies.  Therefore, the purpose of this study is to evaluate the heart rate measurement validity of three

123   consumer photoplethysmographic heart rate monitors compared to an accepted criterion device in

124   accordance with current standards of ANSI/CTI.

125   Methods

126   Participants

127   Twenty healthy subjects (10 males and 10 females) voluntarily completed the study.  Subject

128   characteristics are presented in Table 2.  All subjects were educated on the risks of the procedures and

129   gave informed consent prior to the start of the protocols.  Subjects were recruited verbally from faculty,

130   staff, and students within the university or by e-mail through a local running club.  The study was

131   approved by the Institutional Review Board of Georgia State University.

132   Devices

133   Four heart rate measurement devices, three experimental and one criterion device, were used for this

134   study.  All devices were updated with the most recent software and firmware prior to the start of the

135    study.  No further updates were installed on the devices during the data collection period so that

136    firmware and software remained consistent throughout the study.  The criterion device was the Polar

137    H10 (H10; Firmware 3.0.50, Polar Electro, Kempele, Finland), an electrode chest-strap heart rate

138    monitor.  The Polar H10 uses existing technology from its predecessor Polar H7 which has been

139    validated as above 99% accurate compared to ECG in previous studies (20).  The three PPG experimental

140    devices were the Polar OH 1 (OH1; Firmware 2.0.10, Polar Electro, Kempele, Finland), Apple Watch 4

141    (AW4; Watch OS 5.3.2, Apple, Inc., Cupertino, California) and Garmin Forerunner 945 (FR945; Firmware

142    2.80, Garmin Ltd., Schaffhausen, Switzerland).  The device placement locations were consistent between

143    subjects with OH1 located on the right anterior forearm, AW4 on the right wrist, GF945 on the left wrist.

144    The H10 was fitted on the anterior thorax at the level of the xiphoid process with conduction gel to

145    ensure signal transmission.

146    Procedures

147    Data collection for each subject was completed in a single session and devices were not moved from

148    their specific placement location throughout the entirety of the session.  Subjects arrived at the Applied

149    Physiology Laboratory at Georgia State University or the headquarters of a local running club according

150    to their preferred location. After subjects completed informed consent, investigators recorded

151    anthropometric information including subject-described Fitzpatrick score for skin tone, body mass via

152    calibrated digital scale, body fat percentage via 3-site skinfold test, age, and sex.  Subjects were then

153    verbally informed of the study protocol, which was a running and walking protocol completed on

154    Woodway treadmills, a Pro XL at the university laboratory and a Desmo S at the local running club

155    (Woodway USA, Inc., Waukesha, WI).  Subjects reported a general training intensity level (intensity)

156    described as moderate, high, very high, or elite intensity based upon personal preference and abilities.

157    Walking and running intensities were then assigned by investigators based on this information.    Details

158    about the intensity levels are depicted in Table 1.

159     Table 1. Treadmill Intensities

| Testing Condition | Intensity | | | |
|---|---|---|---|---|
| | Moderate | High | Very High | Elite |
| 2 - Steady State Walk | 2.5 MPH | 2.7 MPH | 3.0 MPH | 3.3 MPH |
| 3 - Steady State Run | 5.0 MPH | 6 MPH | 7.0 MPH | 8.0 MPH |
| 4 - Dynamic (Run/Faster/Fastest) | 5/5.5/6.0 MPH | 6/6.7/7.3 MPH | 7/7.7/8.3 MPH | 8/8.7/9.3 MPH |

160

161     Testing Conditions

162     For each subject, data collection was completed for all 4 testing conditions in a single session. Each

163     testing condition included 1 minute of quiet sitting both prior to and after the treadmill protocol.  SED,

164     WLK, and RUN were 7 minutes in length, including the quiet sitting.  DYN was 12 minutes in length,

165     including the quiet sitting.  For SED, subjects remain seated and motionless for 5 minutes.  For the WLK

166     and RUN, subjects completed 5 minutes of activity at the assigned treadmill speed intensity, which

167     investigators set manually for each trial.  For DYN, a time-based running and walking protocol, each

168     treadmill was identically pre-programmed with 4 different programs to adjust speed at specific time

169     intervals according to the assigned intensity as seen in Table 1.  Walking speed during DYN matched the

170     same intensity speed as WLK condition.  Between each testing condition, subjects rested for 5 minutes

171     to allow heart rate to return to normal.

172     Data Acquisition

173     Data from the OH1 and H10 were transmitted from the device via Bluetooth to an iPad Mini running the

174     Performtek app (Valencell, Inc. Raleigh, North Carolina).  The Performtek app allows for connection of

175     multiple devices and records device data, including heart rate, for side-by-side comparison.  Data from

176    the AW4 was downloaded to RunGap software (CTRL-N ApS, Skødstrup, Denmark) which was then

177    converted to .csv format and imported to Excel. The AW 4 could not be adjusted to record at a specific

178    frequency and required manual data alignment with the same time points of the criterion device for

179    proper analysis.   The GF945 data were downloaded as a raw data file (.tcx file) via device sync with

180    Garmin Connect.  The H10, OH1, and GF 945 were all programmed to record heart rate at 1 Hz. OH1,

181    H10, and GF945 data were then converted to .csv and imported into a Microsoft Excel (Microsoft

182    Corporation, Redmond Washington) spreadsheet for analysis.

183    Statistical Analysis

184    After being organized in Excel, data were imported into SPSS 27 (SPSS; IBM Corporation, Armonk, NY)

185    for further analysis.  Mean Absolute Difference (MAD) and Mean Absolute Percent Error (MAPE) were

186    calculated for each device for each protocol in Excel.  T-tests for the difference between experimental

187    device and criterion device for each stage of each protocol were conducted in SPSS to determine mean

188    difference and standard deviation.  Pearson's R correlation and intraclass correlation (ICC) were

189    calculated to determine general correlation between devices and absolute agreement between devices,

190    respectively.  Lastly, Bland-Altman plots were created with mean bias and upper and lower limits of

191    agreement.  ANSI/CTA standards deem any device with a MAPE ≤ 10% as valid.

192    Results

193    Subject Characteristics

194    Basic subject characteristics are presented in Table 2.  Recruitment of the subject population was

195    coordinated to adhere to the ANSI/CTA standards for device research such that the minimum

196    percentages of subjects met criteria for Body Mass Index (BMI), Fitzpatrick Score (i.e., skin tone), and

197    sex.  The standards as of the ANSI/CTI-2065 were (over the age of 18), sex (no less than 40%

198    male/female), skin tone (minimum 25% from lighter scale and minimum 25% from darker scale), and

9

199    body mass (minimum 10% above 25 kg/m$^2$ and minimum 10% below 20 kg/m$^2$).  Additionally, a

200    minimum of 20 subjects is recommended.

201    Table 2.  Subject characteristics

| Sex | BMI (kg/m^2) | Fitzpatrick Score 1-3 (n) | Fitzpatrick Score 4-6 (n) | Height (m) | Weight (kg) | Body Fat (%) |
|---|---|---|---|---|---|---|
| Male (n = 10) | 24.86 | 6 | 4 | 1.73 | 74.54 | 13.23 |
| Female (n = 10) | 23.72 | 8 | 2 | 1.64 | 63.62 | 28.01 |
| All Subjects | 24.3 | 14 | 6 | 1.7 | 69.1 | 20.6 |

202

## General Device Results

204    Results for all devices can be seen in Tables 3 and 4.  More detailed device results based on specific test

205    conditions can be seen in Appendix 1. Both MAD and MAPE are device HR to criterion HR comparisons

206    for all subjects during the entire 7 or 12 minutes of each testing condition.  The 7-minute testing

207    conditions had approximately 420 data points (HR measurements) per subject and the 12-minute testing

208    conditions had approximately 720 data points per subject.  As the AW4 did not allow for 1Hz HR

209    recording, data points were fewer resulting in approximately 220 data points per subject for the 7-

210    minute protocols and approximately 365 data points per subject for the 12-minute protocol.  MAPE

211    must be ≤ 10% to be considered valid according to ANSI/CTA-2065 standards.   Using this this threshold,

212    each device was considered valid for each condition tested, although the devices did produce differing

213    results for both MAD and MAPE.  Bland-Altman plots for each device's data aggregated across all

214    conditions can be seen for AW4, FR945, and OH1 in Figures 1, 2, and 3, respectively.   Device by test

215    condition Bland-Altman plots can be seen in Appendix 1.

216    Table 3. Mean Absolute Deviation (MAD)

| Protocol | Polar OH1 | Apple Watch 4 | Garmin FR945 |
|---|---|---|---|
| Sedentary | 1.31 | 1.33 | 2.03 |
| Walk | 2.79 | 2.58 | 5.19 |
| Run | 4.00 | 4.29 | 6.51 |
| Dynamic | 2.60 | 2.44 | 7.18 |

217

218 Table 4. Mean Absolute Percent Error (MAPE)

| Protocol | Polar OH1 | Apple Watch 4 | Garmin FR945 |
|---|---|---|---|
| Sedentary | 1.78% | 1.89% | 2.81% |
| Walk | 3.15% | 3.18% | 5.93% |
| Run | 3.43% | 3.51% | 5.25% |
| Dynamic | 2.05% | 1.95% | 5.47% |

219

220 *Polar OH1*

221 The OH1 resulted in a MAD between 1.31 (SED) and 4.00 (RUN) with a MAPE between 1.78% (SED) and

222 3.43% (RUN).  The Bland Altman plot for the OH1 can be seen in Figure 1.  The LoA for the OH1 ranged

223 between -9.406 and 10.586 with a mean bias of .59. The ICC of the OH1 was .991 with 95% CI of .992

224 and .991.

225 **Fig 1. Bland-Altman Plot of all protocols for Polar OH1** Mean bias of 0.59 with upper and lower limits of

226 agreement of 10.586 and -9.406, respectively.

227 *Apple Watch 4*

228 The AW4 produced a MAD between 1.33 (SED) and 4.29 (RUN).  The MAPE for the AW4 ranged between

229 1.89% (SED) and 3.51% (RUN).  The Bland Altman plot for the AW4 for all protocols (SED, WLK, RUN,

230 DYN) can be seen in Figure 2.  The overall Limits of Agreement (LoA) ranged from to -13.314 to 13.974

11

231 with a mean bias of .33. Intraclass Correlation (ICC) was high at .990 with a 95% Confidence Interval (CI)

232 of .990 and .989, upper and lower, respectively. Pearson's r was .990.

233 **Fig 2. Bland-Altman Plot of all protocols for Apple Watch 4** Mean bias of 0.33 with upper and lower

234 limits of agreement of 13.974 and -13.314, respectively.

235 *Garmin Forerunner 945*

236 The FR945 results yielded a MAD between 2.03 (SED) and 7.18 (DYN) with a MAPE range between 2.81%

237 (SED) and 5.93% (WLK). The Bland Altman plot for all protocols can be seen in Figure 3. The FR945

238 produced LoA between -17.269 and 20.469 with a mean bias of 1.6. The ICC was .967 with a 95% CI of

239 .970 and .965. Pearson's r was .969.

240 **Fig 3. Bland-Altman Plot of all protocols for Garmin Forerunner 945** Mean bias of 1.60 with upper and

241 lower limits of agreement of 20.469 and -17.269, respectively.

242 ## Discussion

243 ### ANSI/CTA Standards Validity

244 The principal findings of our study were that HR measurement via PPG technology in the Polar OH1,

245 Apple Watch 4, and Garmin Forerunner 945 met the criteria to be considered valid by the ANSI/CTA

246 standards. All three devices had a MAPE <10% while being evaluated across a broad subject group

247 comprised of adequate representation across various skin tones, BMI levels, and sex.

248 Over the past two decades, wearable fitness devices have progressed in both use and functionality

249 resulting in a broad range of options for consumers. A major advancement is the integration of PPG

250 technology into the devices. By establishing the proper color of the light-emitting diode and refining

251 proprietary algorithms, manufacturers can now provide end-users with myriad physiological information

252    in a single device without the need for a chest strap.  The devices used in this study all use similar PPG

253    technology, primarily differing in only the number of diodes and the manufacturer's unique algorithms.

254    The development of the ANSI/CTA standards for determining device validity defines a framework that

255    generally allows for a more equitable and diverse application of the device characteristics to the total

256    population.  This study represents one of the first studies that has developed the study design in strict

257    accordance with the ANSI/CTA standards.  Subject selection was not random, but instead, individuals

258    were specifically recruited to meet the minimum percent of subject group standards such that age, sex,

259    skin tone and BMI were all adequately represented in the subject group.  Additionally, exercise

260    conditions were specifically designed to adhere to the standards, and subject input was utilized to

261    appropriately set intensities across a very diverse group of subjects.  Although specific analysis of

262    appropriate intensity matching is beyond the scope of this research, visual analysis of the data suggests

263    that all subjects performed each test condition in alignment with the information provided.  Therefore,

264    by implementing a strict study design and appropriately selecting subjects based on the prescribed

265    framework, the results of this study can be broadly applied to the general population.

266    Comparison with Previous Studies

267    Previous studies have attempted to determine validity for various PPG devices, although to our

268    knowledge, this is the first study to strictly apply the ANSI/CTA standards to subject selection and study

269    design.  The devices in this study have been directly and indirectly studied in conjunction with other

270    devices or using different methodologies.  As the consumer electronics market is constantly progressing

271    and new devices are introduced to consumers fairly frequently, direct device comparison is limited and

272    requires inclusion of different versions or generations of the devices.  Although device manufacturers

273    have been researched extensively during the past 6 to 7 years, precise comparison between this study's

274    devices and previous research is very limited.  Of the devices tested in this study, the OH1 has had been

275     researched the most.  This is most likely because the OH1 has stayed consistent during its lifetime

276     whereas other products have had generational changes or complete updates to the product line.  The

277     original OH1 was released in 2017 with only one major upgrade to the OH1 Plus (allowed ANT+

278     communication).  At the same time, Apple has released 4 different Apple Watches, and Garmin

279     progressively released new watches in the Forerunner series with the FR945 being released in late 2019.

280     *Polar OH1*

281     Multiple studies have previously provided ample evidence of the validity of the OH1.  Schubert et al.

282     found the mean bias to be slightly higher than the current study (.59 versus .76) but a narrower LoA

283     (-9.406 and 10.586  versus -3.83 to 5.35), but is limited in application as the study compared only a

284     mean heart rate for a yoga session while also suffering from unbalanced subject sex selection (n=15, 3

285     males) with limited BMI and Fitzpatrick Scale variation (3) A more recent study found a lower mean bias

286     (.27) and narrower LoA (-4.68 to 5.22) than the current study (4).  Direct comparison is difficult as

287     subjects the previous study noted all subjects held the handrail potentially decreasing any motion

288     artifact, and the study also lacked any diversity with Fitzpatrick Scale and BMI.  A 2019 study assessing

289     different activities resulted in lower biases for walking and running (.18 versus .41 and .37 versus 1.28,

290     respectively) but this study was biased towards males (n=70, males = 54), did not report BMI, and

291     although it referenced skin tone, specific subject representation of skin tone levels was not reported (6).

292     Additionally, A more recent study has further confirmed the validity of the OH1 in various activities, and

293     across all activities found a higher mean bias (1 versus .59), a broader LoA (-20 to 19 versus -9.406 to

294     10.586) with a lower *r* (.957 versus .991)  compared to the current study, but like other studies lacked

295     subject information about skin tone and BMI (13).

296     *Apple Watch 4*

297     Apple regularly releases new products on an annual basis.  As such, direct assessment of the AW4 is

298     difficult, but evaluation of previous versions is available in the research.  Dooley at al. evaluated the

14

299     first-generation Apple Watch across a wide range of BMI and exercise intensities finding a higher MAPE

300     for walking (5.60% versus 3.18%) and running (6.70% versus 3.51%) compared to the current study (7).

301     Although the study utilized different treadmill walking and running intensities, the heart rate data was

302     only recorded for a single time point and Fitzpatrick Scale was not recorded.  In 2019 Hwang et al.

303     researched the Apple Watch 2, revealing a much tighter LoA (-6.0 to 3.9 versus -13.314 to 13.974) but a

304     slightly lower mean bias (-1.0 versus .330) than the current study, but Hwang used a model of

305     tachycardia with electrical pacing so direct comparisons are difficult to determine (8).  Nelson et al.

306     released research concerning the Apple Watch 3 in 2019 (9), finding a higher mean bias (1.80 versus

307     .330) and higher MAPE (5.86% versus 2.63%) than our study but Nelson's study was a single-subject

308     free-living design comparing different devices.  Lastly, Duking et al. investigated the AW4 but the

309     authors did not calculate key validity metrics (MAD, MAPE, ICC) with only a slightly lower $r$ (.97 versus

310     .984) available for comparison to the current study  (10). Only one of these studies actively recruited

311     subjects with skin tone variations, but the delineation was limited to white and non-white and

312     ethnicity/race, not a skin tone scale (7).

313     *Garmin Forerunner 945*

314     FR945 validity data is lacking in the literature.  The prior device-specific research that is available

315     generally concerns the Forerunner 235 versus this study's 945.  In Dooley's 2017 study, the Garmin

316     Forerunner 235 had large deviations from the criterion HR with as high as MAPE of 24.38% (2.81% to

317     5.93% for the current study).  In 2019, Stove et al. also completed research on Garmin Forerunner 235

318     validity revealing much lower ICC values (.480 to .905 versus .895 to .973) compared to the current

319     study, but had a limited number of heart rate data points as data was only recorded once per minute

320     (11).

### Device Differences

322    Although all devices tested were deemed valid, differences in MAD and MAPE for different devices did

323    exist.  Whether these differences are functionally important is determined by the consumer.  In respect

324    to the criterion, the OH1 and AW4 tended to have lower differences for MAD and MAPE values, as well

325    as a higher ICC and narrower CI compared to the FR945.  Similarly, the ICC and r were lower for the

326    FR945 compared to the other two devices with the OH1 having a very slightly higher ICC and r than the

327    AW4.

328    The reasons for these differences could be due to multiple factors.  First, as previously mentioned, the

329    devices all differ for functionality and intended use.  Secondly, although each device was worn according

330    to manufacturer's specifications, devices differed in wristband/armband material and the size of the

331    recording device.  The FR945 and AW4 are both wrist-worn monitors but differing styles and materials

332    of the wristbands resulted in slightly different fitment for the devices on individual subjects due to

333    variation in wrist diameter.  The OH1 had the smallest recording device and was secured to the lower

334    arm via a fabric elastic band.  Although the technology for the PPG light-emitting diode, appears to be

335    similar between devices, individual devices variances between the number of diodes and spacing of

336    diodes is visually apparent.  The most likely reason for the differences, though, is the manufacturer-

337    specific algorithm that converts the PPG raw data to heart rate information.   Other differences in

338    proprietary technology, such as the device specific hardware and software for recording and processing

339    also presumably exist.

### Limitations

341    Although this study was conducted according to the current ANSI/CTA standards, certain limitations do

342    exist.  First, the subject group (n=20) is considered the minimum subject group size and minimum

343    percentage for the specific parameters of BMI and Fitzpatrick Scale.  Although adequate for ANSI/CTA

16

344      standards, future studies should consider a larger subject group so that those two parameters can be

345      more intricately analyzed within the total subject group. Additionally, the results of this study can only

346      be applied to the specific devices and their corporation-specific algorithms to compute heart rate from

347      PPG signals. As the technology continues to advance, it is plausible that the corporations will refine the

348      algorithms in attempts to improve validity. Lastly, the ANSI/CTA standards place limitations on the

349      subject group such that individuals with tattoos in the sensor location should not be included in the

350      study due to presumed alterations in how the photosensor reads the reflection of the capillary beds. As

351      it can be argued that tattoos on the arm and wrist have become popularized as of late, the validity of

352      these devices cannot be confirmed in this subgroup.

353      ## Conclusions

354      As consumers are consistently utilizing a variety of devices to track health metrics which rely on heart

355      rate measurements, it is vital that the PPG recording technology and manufacturer proprietary

356      algorithms properly represent the actual heart rate of the individual. As the end-consumer of these

357      devices represents a wide range of subject characteristics, it is equally important that the devices

358      correctly record heart rate across variations in age, skin tone, sex, and BMI. By utilizing the ANSI/CTA

359      standards for heart rate recording devices, this and future studies can be more confident that the data

360      recorded by the device can be utilized confidently by the majority of the population. In this study, one

361      of the first to implement a study design in accordance with the ANSI/CTA standards, the Polar OH1,

362      Apple Watch 4, and Garmin Forerunner 945 were all deemed valid in their measurement of heart rate.

363      Consumers of various age, sex, body composition, and skin tone can be confident that the heart rate

364      data presented to them is within a strict range for validity and represents their unique characteristics.

365

## Acknowledgements

The authors would like to acknowledge Phung Tran, Robert Tippett, Jr., Shreya Kulkarni for their

assistance in data collection.  Additionally, the authors would like to acknowledge Greg Dodd for his

assistance with data alignment between the AW4 and H10.  Lastly, the authors would also like to thank

the Dr. David E. Martin Sport Science Research Fund and the Atlanta Track Club for their respective

financial and logistical contributions to the study.

## Conflict of Interest

The authors state no conflict of interests.

## Abbreviations

ANSI/CTA: American National Standards Institute/Consumer Technology Association

AW4: Apple Watch 4

BPM: Beats per minute

GF945: Garmin Forerunner Multi-function watch

H10: Polar H10 Chest Strap Heart Rate Monitor (criterion)

ICC: Intraclass correlation

MAD: Mean Absolute Deviation

MAPE: Mean Absolute Percent Error

OH1: Polar OH1 Armband Heart Rate Monitor

PPG: Photoplethysmography

388   References

389   1.      Thompson WR. WORLDWIDE SURVEY OF FITNESS TRENDS FOR 2020. ACSM's Health & Fitness

390   Journal. 2019;23(6):10-8.

391   2.      Allen J. Photoplethysmography and its application in clinical physiological measurement. Physiol

392   Meas. 2007;28(3):R1-39.

393   3.      Schubert MM, Clark A, De La Rosa AB. The Polar (®) OH1 Optical Heart Rate Sensor is Valid

394   during Moderate-Vigorous Exercise. Sports Med Int Open. 2018;2(3):E67-e70.

395   4.      Hettiarachchi IT, Hanoun S, Nahavandi D, Nahavandi S. Validation of Polar OH1 optical heart rate

396   sensor for moderate and high intensity physical activities. PLoS ONE. 2019;14(5):1-13.

397   5.      Olstad BH, Zinner C. Validation of the Polar OH1 and M600 optical heart rate sensors during

398   front crawl swim training. PLoS One. 2020;15(4):e0231522.

399   6.      Hermand E, Cassirame J, Ennequin G, Hue O. Validation of a Photoplethysmographic Heart Rate

400   Monitor: Polar OH1. Int J Sports Med. 2019;40(7):462-7.

401   7.      Dooley EE, Golaszewski NM, Bartholomew JB. Estimating Accuracy at Exercise Intensities: A

402   Comparative Study of Self-Monitoring Heart Rate and Physical Activity Wearable Devices. JMIR Mhealth

403   Uhealth. 2017;5(3):e34.

404   8.      Hwang J, Kim J, Choi KJ, Cho MS, Nam GB, Kim YH. Assessing Accuracy of Wrist-Worn Wearable

405   Devices in Measurement of Paroxysmal Supraventricular Tachycardia Heart Rate. Korean Circ J.

406   2019;49(5):437-45.

407   9.      Nelson BW, Allen NB. Accuracy of Consumer Wearable Heart Rate Measurement During an

408   Ecologically Valid 24-Hour Period: Intraindividual Validation Study. JMIR Mhealth Uhealth.

409   2019;7(3):e10828.

410    10.    Düking P, Giessing L, Frenkel MO, Koehler K, Holmberg HC, Sperlich B. Wrist-Worn Wearables for

411    Monitoring Heart Rate and Energy Expenditure While Sitting or Performing Light-to-Vigorous Physical

412    Activity: Validation Study. JMIR Mhealth Uhealth. 2020;8(5):e16716.

413    11.    Støve MP, Haucke E, Nymann ML, Sigurdsson T, Larsen BT. Accuracy of the wearable activity

414    tracker Garmin Forerunner 235 for the assessment of heart rate during rest and activity. J Sports Sci.

415    2019;37(8):895-901.

416    12.    Gillinov S, Etiwy M, Wang R, Blackburn G, Phelan D, Gillinov AM, et al. Variable Accuracy of

417    Wearable Heart Rate Monitors during Aerobic Exercise. Med Sci Sports Exerc. 2017;49(8):1697-703.

418    13.    Muggeridge DJ, Hickson K, Davies AV, Giggins OM, Megson IL, Gorely T, et al. Measurement of

419    Heart Rate Using the Polar OH1 and Fitbit Charge 3 Wearable Devices in Healthy Adults During Light,

420    Moderate, Vigorous, and Sprint-Based Exercise: Validation Study. JMIR Mhealth Uhealth.

421    2021;9(3):e25313.

422    14.    Chow HW, Yang CC. Accuracy of Optical Heart Rate Sensing Technology in Wearable Fitness

423    Trackers for Young and Older Adults: Validation and Comparison Study. JMIR Mhealth Uhealth.

424    2020;8(4):e14707.

425    15.    Baek S, Ha Y, Park HW. Accuracy of Wearable Devices for Measuring Heart Rate During

426    Conventional and Nordic Walking. Pm r. 2020.

427    16.    Preejith SP, Annamol A, Joseph J, Sivaprakasam M. Design, development and clinical validation

428    of a wrist-based optical heart rate monitor. 2016 IEEE International Symposium on Medical

429    Measurements and Applications (MeMeA). 2016:1-6.

430    17.    Boonya-Ananta T, Rodriguez AJ, Ajmal A, Du Le VN, Hansen AK, Hutcheson JD, et al. Synthetic

431    photoplethysmography (PPG) of the radial artery through parallelized Monte Carlo and its correlation to

432    body mass index (BMI). Scientific reports. 2021;11(1):2570.

433    18.    Sañudo B, De Hoyo M, Muñoz-López A, Perry J, Abt G. Pilot Study Assessing the Influence of Skin

434    Type on the Heart Rate Measurements Obtained by Photoplethysmography with the Apple Watch. J

435    Med Syst. 2019;43(7):195.

436    19.    Addison PS, Jacquel D, Foo DMH, Borg UR. Video-based heart rate monitoring across a range of

437    skin pigmentations during an acute hypoxic challenge. J Clin Monit Comput. 2018;32(5):871-80.

438    20.    Cheatham SW, Kolber MJ, Ernst MP. Concurrent validity of resting pulse-rate measurements: a

439    comparison of 2 smartphone applications, the polar H7 belt monitor, and a pulse oximeter with

440    bluetooth. Journal of sport rehabilitation. 2015;24(2):171-8.

441

442    Appendix 1.

443    Polar OH1 Detailed Results

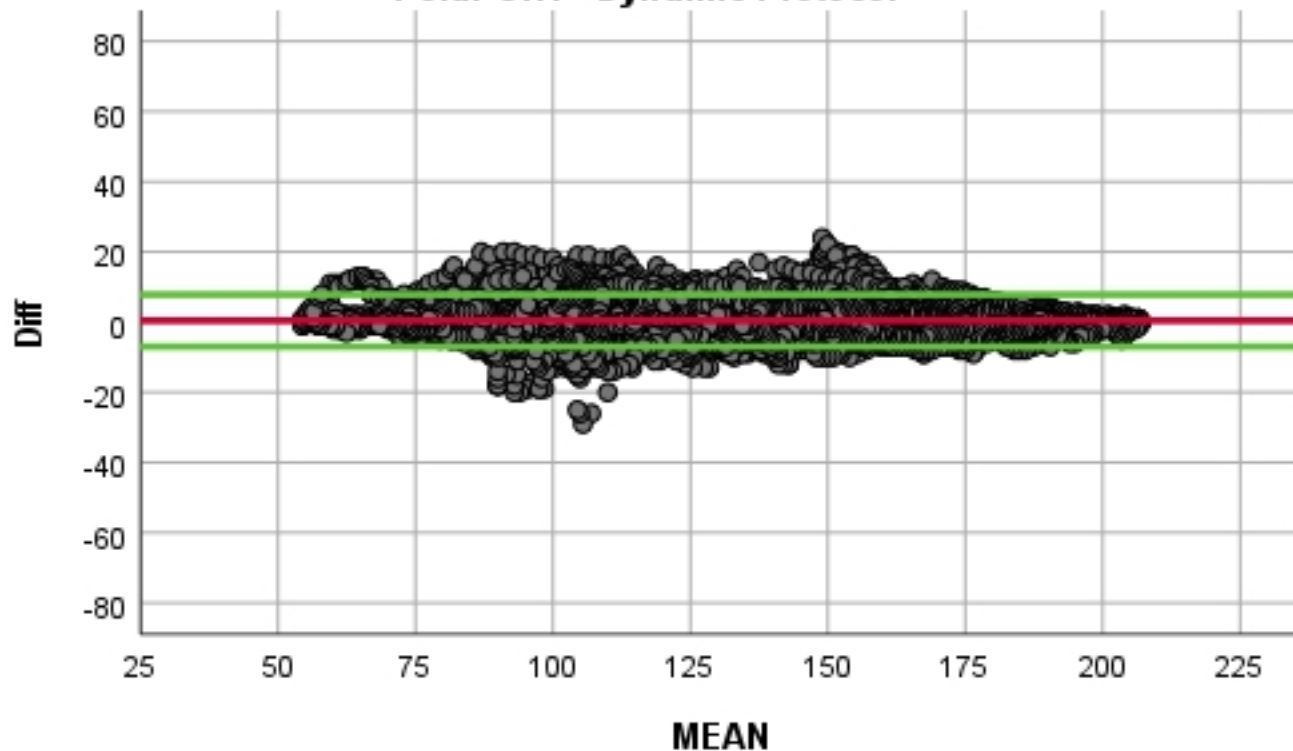| Polar OH1 | SED | WALK | RUN | DYN | ALL |
|---|---|---|---|---|---|
| MEAN | 0.290 | 0.410 | 1.280 | 0.470 | 0.590 |
| SD | 2.015 | 4.817 | 8.398 | 3.778 | 5.100 |
| LoA Upper | 4.239 | 9.851 | 17.740 | 7.875 | 10.586 |
| LoA Lower | -3.659 | -9.031 | -15.180 | -6.935 | -9.406 |
| LoA Range | 7.899 | 18.883 | 32.920 | 14.810 | 19.992 |
|  |  |  |  |  |  |
| ICC | 0.989 | 0.957 | 0.965 | 0.993 | 0.991 |
| 95% Confidence Lower | 0.989 | 0.955 | 0.962 | 0.993 | 0.991 |
| 95% Confidence Upper | 0.990 | 0.959 | 0.967 | 0.993 | 0.992 |
|  |  |  |  |  |  |
| Pearson's R | 0.990 | 0.958 | 0.966 | 0.993 | 0.991 |

444

**Fig 4.  Bland-Altman Plot of Polar OH1 Sedentary Protocol.**  Mean bias of 2.015 with upper and

lower limits of agreement of 4.239 and -3.659, respectively.

**Fig 5.  Bland-Altman Plot of Polar OH1 Walking Protocol.** Mean bias of 4.817 with upper and

lower limits of agreement of 9.851 and -9.031, respectively

**Fig 6.  Bland-Altman Plot of Polar OH1 Running Protocol.**  Mean bias of 8.398 with upper and

lower limits of agreement of 17.740 and -15.180 respectively.

**Fig 7.  Bland-Altman Plot of Polar OH1 Dynamic Protocol.**  Mean bias of 3.778 with upper and

lower limits of agreement of 7.875 and -6.935, respectively.

Apple Watch 4 Detailed Results

| Apple Watch 4 | SED | WALK | RUN | DYN | ALL |
|---|---|---|---|---|---|
| MEAN | 0.170 | -0.880 | 1.580 | 0.420 | 0.330 |
| SD | 2.020 | 4.500 | 12.257 | 5.589 | 6.961 |
| LoA Upper | 4.129 | 7.940 | 25.604 | 11.374 | 13.974 |
| LoA Lower | -3.789 | -9.700 | -22.444 | -10.534 | -13.314 |
| LoA Range | 7.918 | 17.640 | 48.047 | 21.909 | 27.287 |
| | | | | | |
| ICC | 0.990 | 0.960 | 0.926 | 0.985 | 0.984 |
| 95% Confidence Lower | 0.989 | 0.956 | 0.921 | 0.984 | 0.984 |
| 95% Confidence Upper | 0.990 | 0.965 | 0.931 | 0.985 | 0.984 |
| | | | | | |
| Pearson's R | 0.990 | 0.962 | 0.927 | 0.985 | 0.984 |

**Fig 8. Bland-Altman Plot of Apple Watch 4 Sedentary Protocol.** Mean bias of 2.020 with upper and lower limits of agreement of 4.129 and -3.789, respectively.

**Fig 9. Bland-Altman Plot of Apple Watch 4 Walking Protocol.** Mean bias of 4.500 with upper and lower limits of agreement of 7.940 and -9.700, respectively.

**Fig 10. Bland-Altman Plot of Apple Watch 4 Running Protocol.** Mean bias of 12.257 with upper and lower limits of agreement of 25.604 and -22.444, respectively.

**Fig 11. Bland-Altman Plot of Apple Watch 4 Dynamic Protocol.** Mean bias of 5.589 with upper and lower limits of agreement of 11.374 and -10.534, respectively.

23

464    Garmin Forerunner 945 Detailed Results

| Garmin FR945 | SED | WALK | RUN | DYN | ALL |
|---|---|---|---|---|---|
| MEAN | 1.020 | -0.710 | 2.740 | 2.620 | 1.600 |
| SD | 3.106 | 7.590 | 11.158 | 11.760 | 9.627 |
| LoA Upper | 7.108 | 14.166 | 24.610 | 25.670 | 20.469 |
| LoA Lower | -5.068 | -15.586 | -19.130 | -20.430 | -17.269 |
| LoA Range | 12.176 | 29.753 | 43.739 | 46.099 | 37.738 |
| | | | | | |
| ICC | 0.973 | 0.895 | 0.934 | 0.922 | 0.967 |
| 95% Confidence Lower | 0.963 | 0.890 | 0.922 | 0.912 | 0.965 |
| 95% Confidence Upper | 0.979 | 0.900 | 0.943 | 0.931 | 0.970 |
| | | | | | |
| Pearson's R | 0.975 | 0.896 | 0.938 | 0.930 | 0.969 |

465

466    **Fig 12.  Bland-Altman Plot of Garmin Forerunner 945 Sedentary Protocol.**  Mean bias of 3.106

467    with upper and lower limits of agreement of 7.108 and -5.068, respectively.

468    **Fig 13.  Bland-Altman Plot of Garmin Forerunner 945 Walking Protocol.**  Mean bias of 7.590

469    with upper and lower limits of agreement of 14.166 and -15.586, respectively.

470    **Fig 14.  Bland-Altman Plot of Garmin Forerunner 945 Running Protocol.**  Mean bias of 11.158

471    with upper and lower limits of agreement of 24.610 and -19.130, respectively.

472    **Fig 15.  Bland-Altman Plot of Garmin Forerunner 945 Dynamic Protocol.**  Mean bias of 11.760

473    with upper and lower limits of agreement of 20.469 and -17.269, respectively.

24

**Polar OH1 - Dynamic Protocol**

BA OH1 DYN

Apple Watch 4 - Sedentary Protocol

BA AW4 Sed

Apple Watch 4 - Walking Protocol

BA AW4

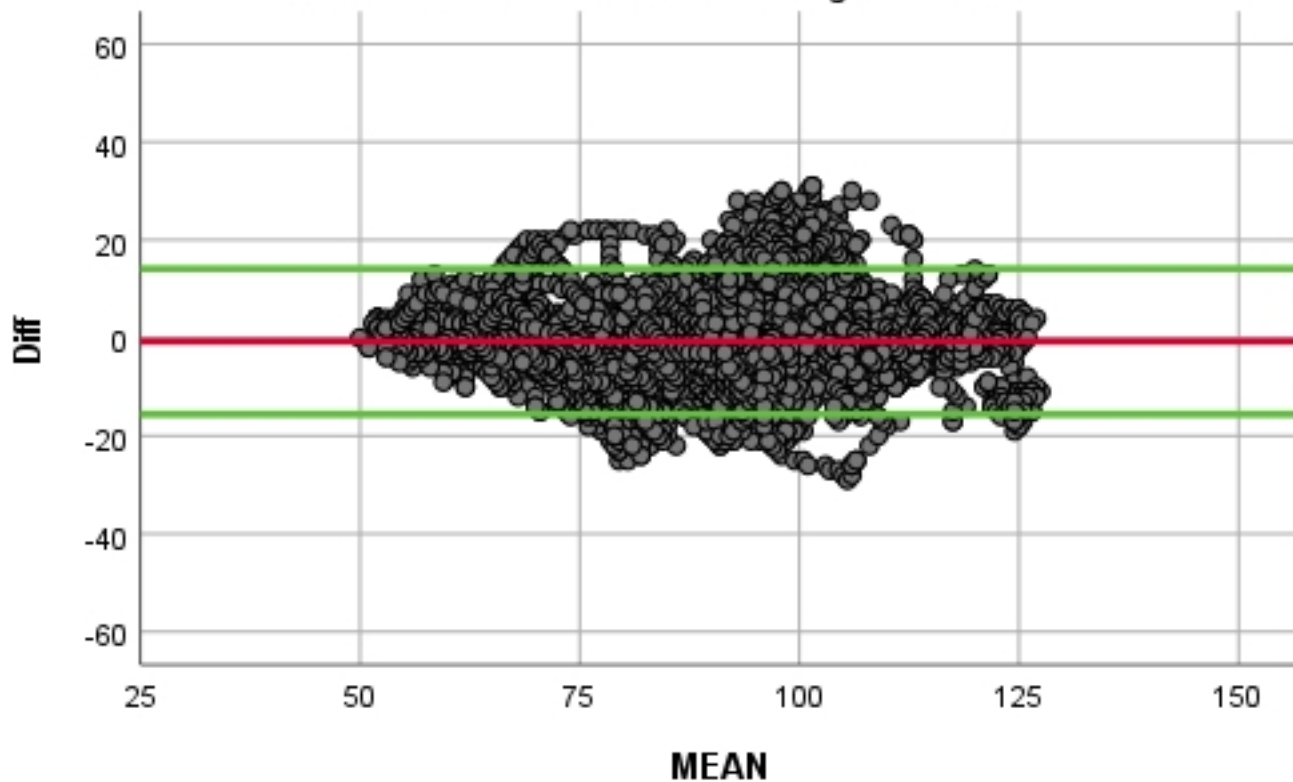Apple Watch 4 - Running Protocol

BA AW4 RUN

**Apple Watch 4 - Dynamic Protocol**
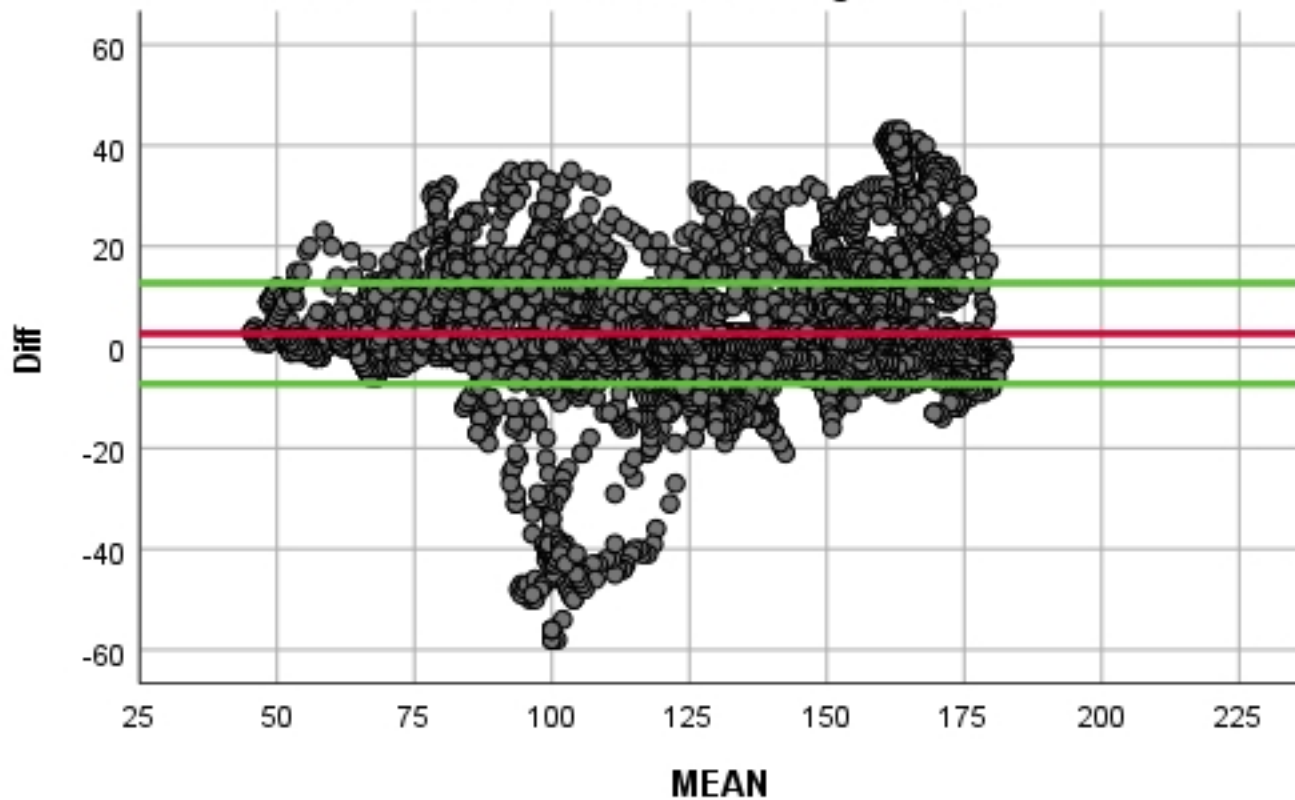
BA AW4 DYN

Garmin Forerunner 945 - Sedentary Protocol
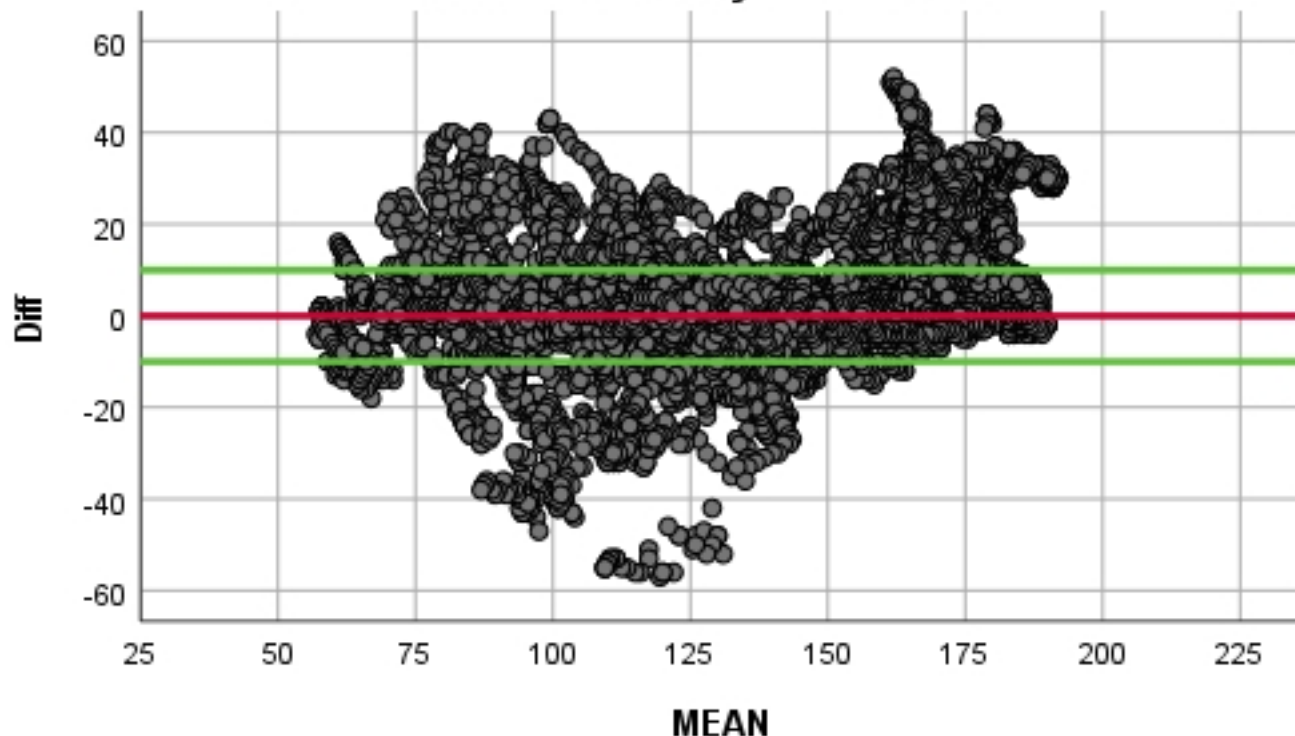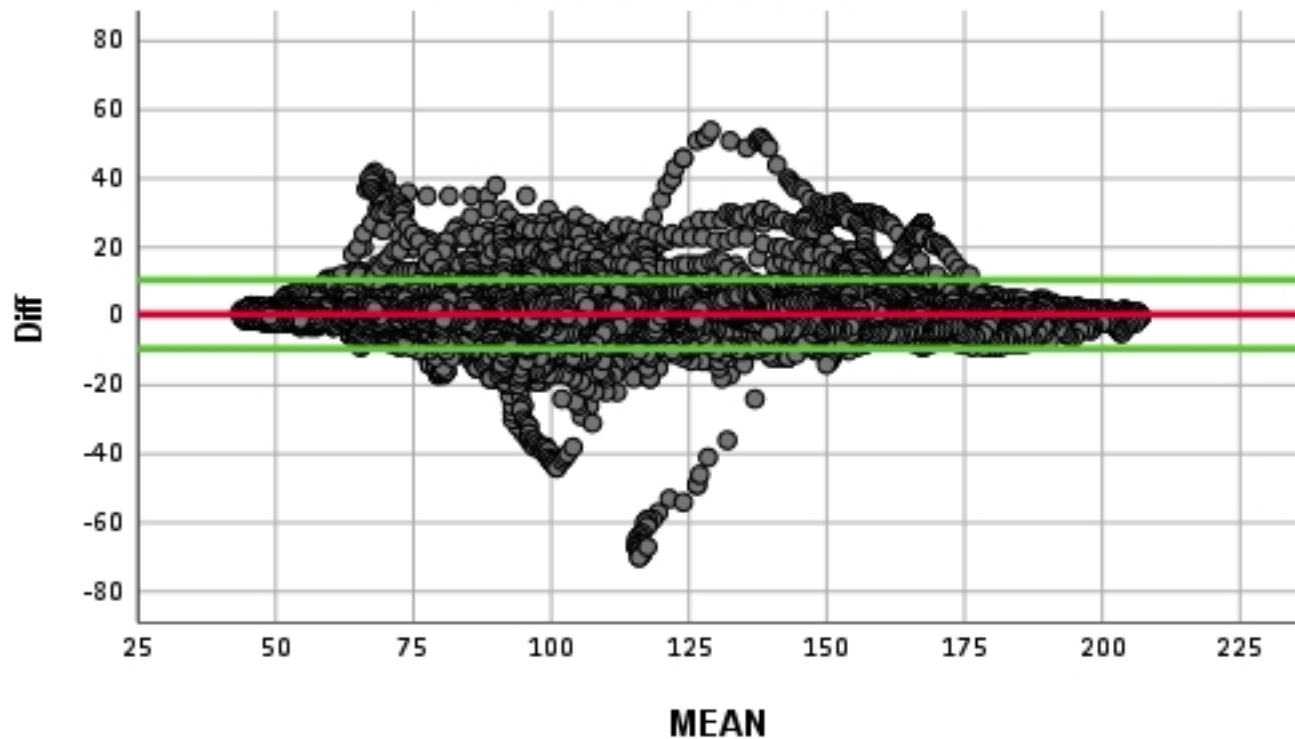
BA FR945 Sed

Gamrin Forerunner 945 - Walking Protocol

BA FR945 WLK

BA FR945 RUN

Garmin Forerunner 945 - Dynamic Protocol
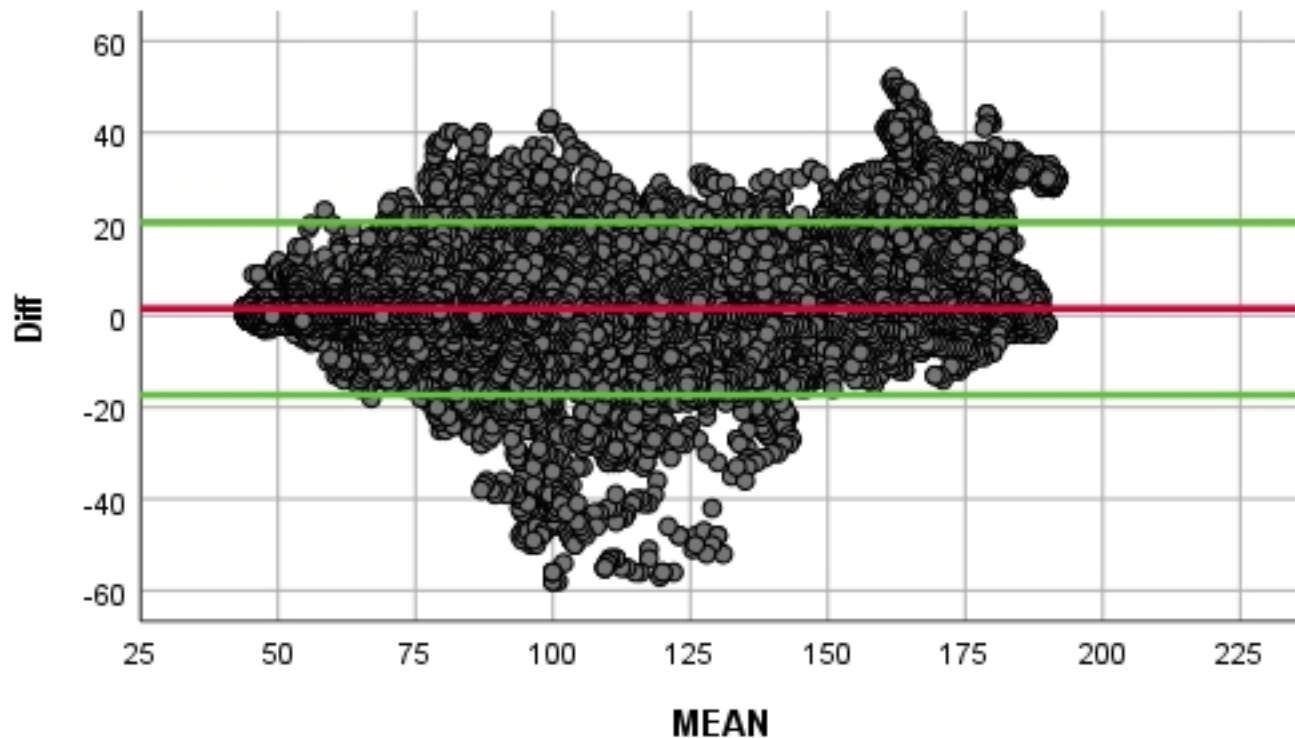
BA FR945 DYN
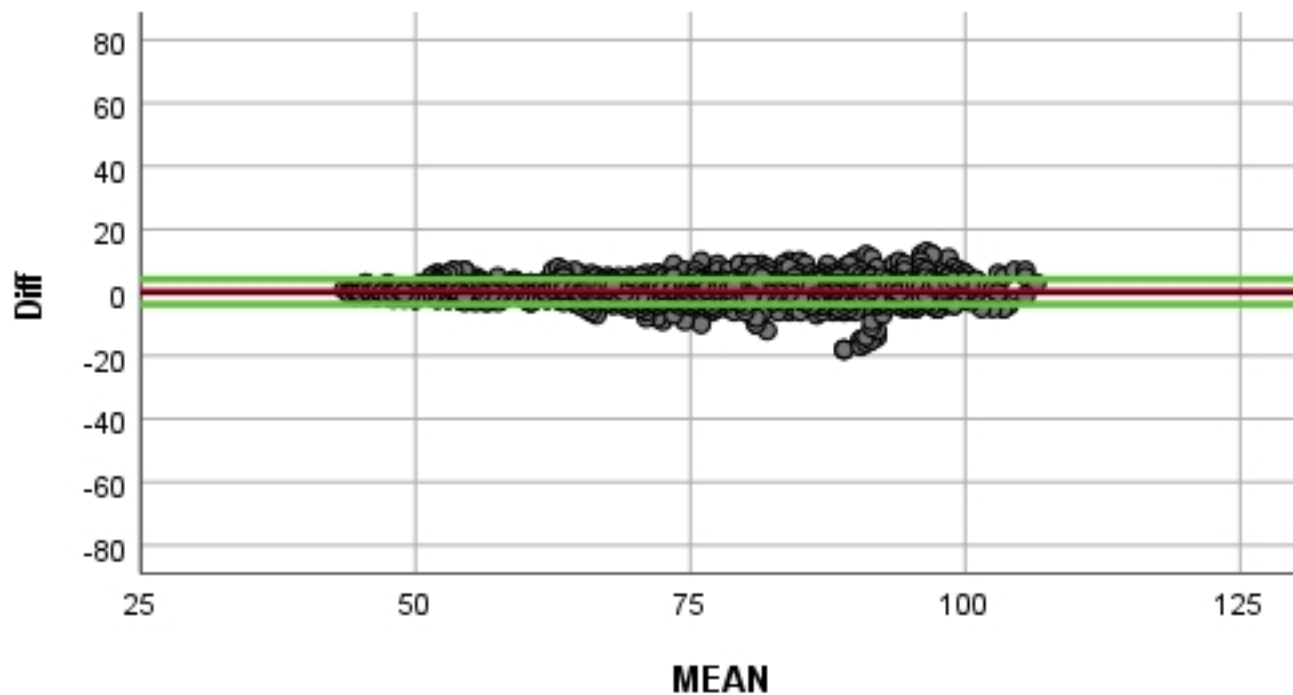
Polar OH1 - All Protocols

BA Plot All OH1
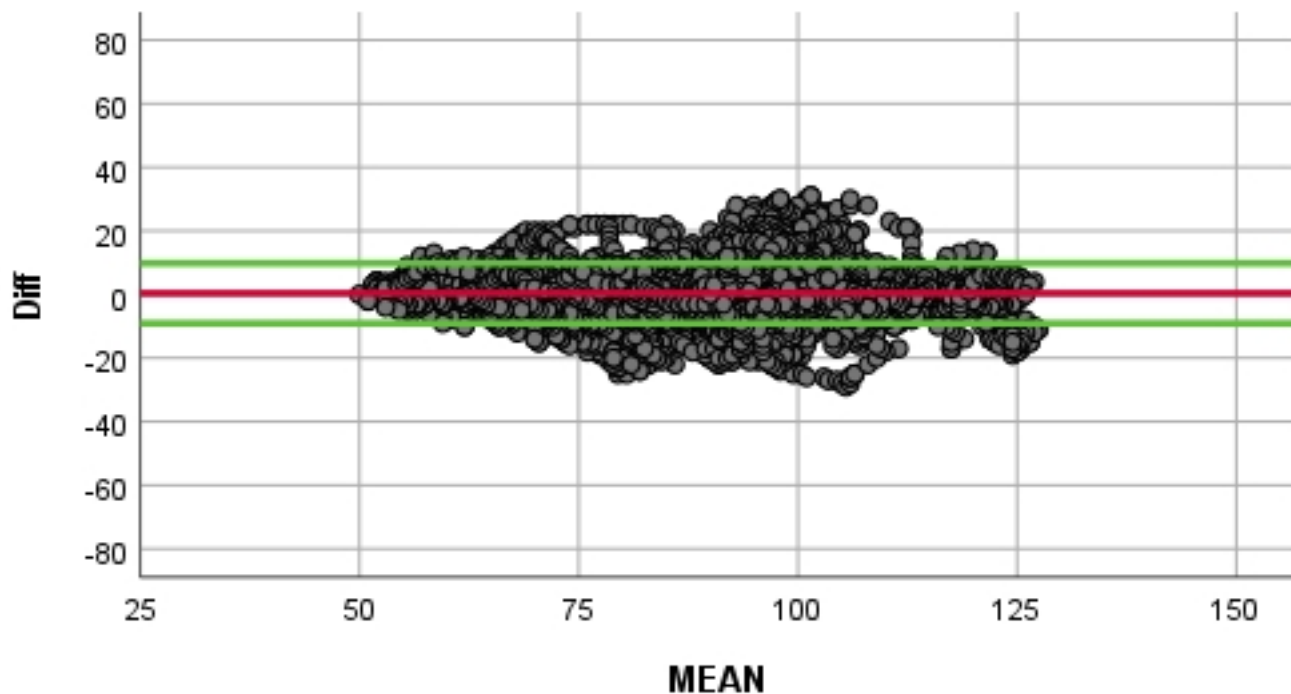
BA Plot All AW4

Garmin Forerunner 945 - All Protocols

BA Plot All FR945
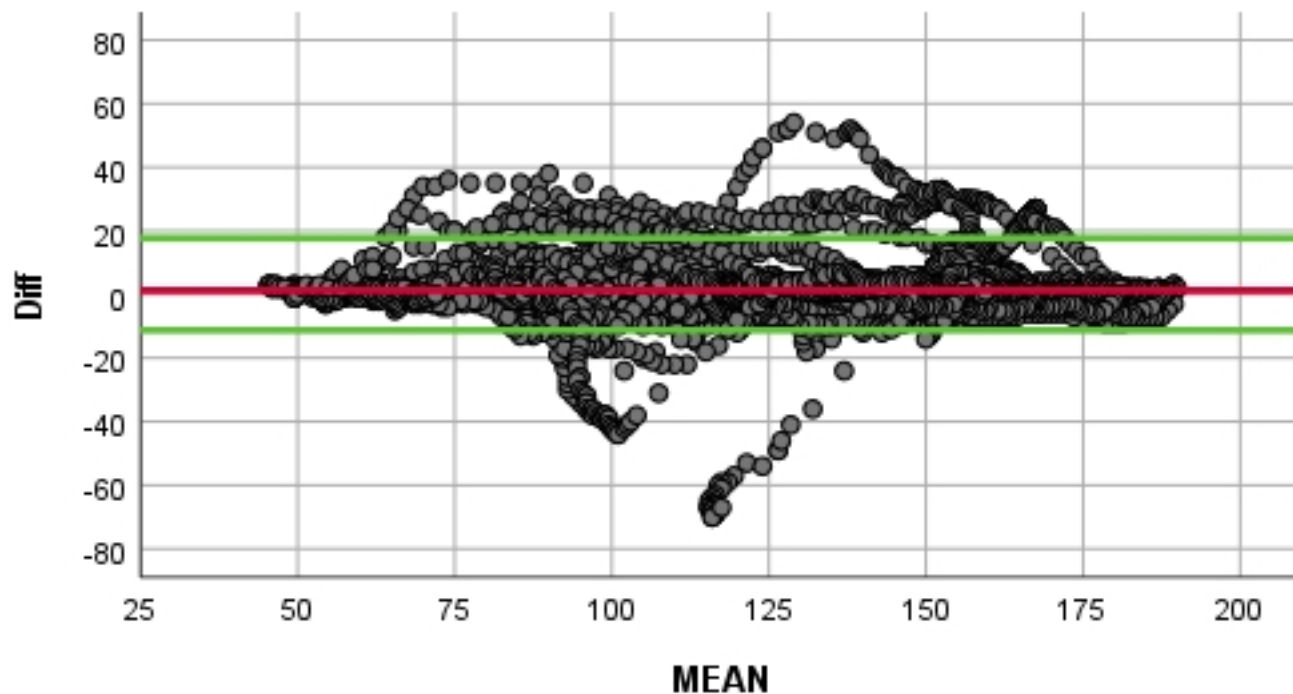
Polar OH 1 - Sedentary Protocol

BA OH1 Sed

Polar OH1 - Walk Protocol

BA OH1 WLK

Polar OH1 - Run Protocol

BA OH1 RUN