

# Validity of the Empatica E4 wristband to estimate resting-state heart rate variability in a lab-based context

**Short version: A lab-based validity assessment of the E4 Empatica**

Stuyck Hans <sup>(1,2)\*</sup>, Dalla Costa Leonardo <sup>(1)</sup>, Cleeremans Axel <sup>(2)</sup> & Van den Bussche Eva <sup>(1)</sup>

<sup>1</sup>*KU Leuven, Faculty of Psychology and Educational Sciences, Brain & Cognition, Tiensestraat 102, 3000 Leuven, Belgium*

<sup>2</sup>*Université libre de Bruxelles, Faculty of Psychology and Education Sciences, Center for Research in Cognition and Neurosciences, Franklin Rooseveltlaan 50, 1050 Brussel, Belgium*

[hans.stuyck@kuleuven.be](mailto:hans.stuyck@kuleuven.be)

[leonardodallacosta.psy@gmail.com](mailto:leonardodallacosta.psy@gmail.com)

[axcleer@ulb.ac.be](mailto:axcleer@ulb.ac.be)

[eva.vandenbussche@kuleuven.be](mailto:eva.vandenbussche@kuleuven.be)

## ABSTRACT

Lab research might benefit from the advantages of wearable devices, such as their ease of use, to estimate pulse rate (PR) and pulse rate variability (PRV) as an equivalent for heart rate (HR) and heart rate variability. However, before implementing them in a lab context, the validity of the PR and PRV, also on ultra-short time scales (e.g., 30s), needs to be confirmed. We recorded heart activity simultaneously with an E4 wristband and an ECG device in a seated resting condition for 5min. Our results showed that HR, RMSSD, SDNN and LF, but not HF, were validly estimated by the E4 wristband. Furthermore, the E4 wristband could validly estimate PR with recording lengths as short as 10s. RMSSD and SDNN were validly estimated using 30s or 120s or an average of multiple short intervals (10s), while HF likely requires longer recording intervals. Based on this study, we formulated several recommendations for using the E4 wristband in a lab context.

**Keywords:** Empatica E4, validity, ultra-short intervals, ECG, heart rate, heart rate variability

## 1. INTRODUCTION

The study of heart rate variability (HRV) has been the object of scientific scrutiny for centuries (for a historical overview, see Berntson et al., 1997). HRV relates to the variation in beat-to-beat intervals (i.e., interbeat interval; IBI). These fluctuations of the IBI result from the active interplay between the parasympathetic (i.e., PNS; rest and digest) and the sympathetic (i.e., SNS; fight or flight) nervous system (Shaffer et al., 2014). In general, it has been found that higher HRV is associated with better adaptability to the environment, for example, concerning physical health, stress regulation, and self-regulation (Graham et al., 2019; Ottaviani et al., 2018; Pulpulos et al., 2018). The most commonly used HRV metrics are situated in the time and frequency domain (Shaffer & Ginsberg, 2017). Time-domain metrics represent the variability of the IBI. For example, the standard deviation of normal IBIs (SDNN) and the root mean square of successive differences between normal IBIs (RMSSD). Frequency domain metrics are represented by the absolute (relative or normalized) power of a signal in one of the four frequency bands related to HRV (i.e., ultra low frequency [ULF;  $\leq 0.003$  Hz], very low frequency [VLF; 0.0033-0.04 Hz], low frequency [LF; 0.04-0.15 Hz] and high frequency [HF; 0.15-0.4 Hz]; Shaffer & Ginsberg, 2017). The informative value of HRV for the individual's mental and physical adaptability has given rise to its use in research settings (see Ottaviani et al., 2018; Stone et al., 2021). HRV is typically measured in laboratory settings by connecting participants to an electrocardiogram device (i.e., ECG). The ECG represents the heart's bio-electrical activity through several waves, of which the large r-peak is used to determine the IBI (Silverthorn, 2004). The ECG's fine-grained representation of the heart's activity makes it the current gold standard to measure HRV (Berntson et al., 1997). However, an ECG device is relatively expensive. Hence, labs often only have one ECG device so that just one participant at a time can be tested. This increases the cost of using ECG in terms of labor cost and meeting sample size requirements. Although ECG studies have been performed on diverse populations (e.g., Ajayi et al., 2021; Nuske et al., 2021), we contend that some targeted populations in lab-based studies might benefit from easier, less time-consuming, and less invasive heart recording setups. For example, young children, older adults, or patients who are already regularly submitted to testing might benefit from a simpler, less obtrusive setup.

One promising solution could be the use of wearable devices in the lab (e.g., Polar smartwatches). These wearables have been developed for daily use, informing users about the complex interaction between behavior, environment, and physical health (Castaneda et al., 2018; Ishaque et al., 2021). These wearable devices are considerably cheaper than an ECG device so that multiple wearables can be bought for the price of one ECG device. Moreover, their design as a watch makes them an almost unnoticeable device that does not hamper its users (Castaneda et al., 2018). However, the estimation of HRV by such devices is quite different from its ECG-based estimation. It is based on the changes in blood volume measured in the body's periphery (e.g., wrist). The differences in blood volume result from the heart contracting to push blood out (i.e., systole) and from the heart's subsequent relaxation (i.e., diastole; Shaffer et al., 2014). These differences in blood volume are measured via photoplethysmography (PPG). This technique measures the amount of light absorbed by the blood vessels, which is proportional to the variations in blood volume. These blood volume variations are represented by the systolic and diastolic pulse waves jointly representing the blood-volume-pulse signal (i.e., BVP; see Figure 1; Alqaraawi et al., 2016). The time interval between the fiducial point on each successive systolic pulse wave is used to estimate the pulse-to-pulse interval (i.e., PPI) on which pulse rate (PR; i.e., beats per minute) and pulse rate variability (i.e., PRV) are calculated, as an equivalent to assessing heart rate (HR) and HRV based on the r-peaks of the ECG (Alqaraawi et al., 2016).

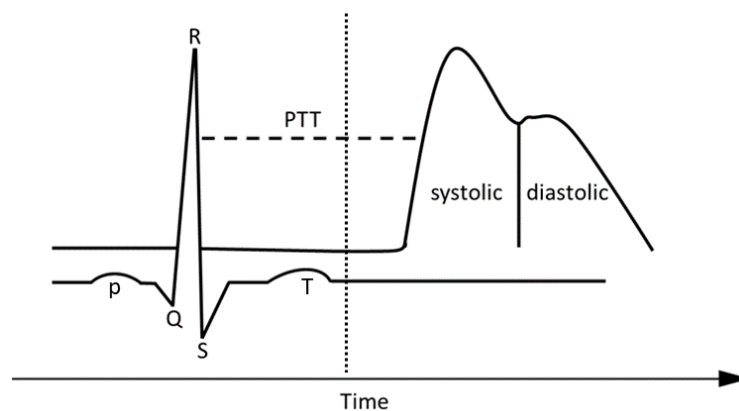
However, there exists a time lag between the depolarization of the ventricles (r-peak) initiating the mass ventricular contraction of the heart and the increase of blood volume measured in the body's

periphery (i.e., the pulse transit time). Features of the blood vessels influence this time lag (e.g., elasticity, vascular diameter) that can change the shape of the pulse waves, resulting in the misalignment between the r-peak and the systolic pulse wave. This might affect the temporal dynamics of the PPI and, thus, the accuracy of the HRV estimate since alterations of the PPI dynamics would not alter PRV but rather the difference between PRV and HRV (Lu & Yang, 2009). Furthermore, several other factors might introduce errors in the accuracy of the PRV measurement. For example, the sampling frequency of such wearable devices is considerably lower than an ECG device (e.g., Empatica E4 wristband = 64Hz). This lower sampling frequency provides a less fine-grained temporal resolution representing both pulse waves (Laborde et al., 2017). A clear and distinct pulse wave morphology is critical to accurately detect the fiducial point on the systolic wave. Also, the PPG-based PRV estimation is more vulnerable to orthostatic changes, motion artifacts, stress induction, and respiratory patterns than ECG-based HRV (see Mejia-Mejia et al., 2020 for a review).

So before implementing such wearable devices in lab research, it is vital to assess how valid the PRV metrics derived by such devices are (Schuurmans et al., 2020). The most common way to assess the validity of PPG-based PR and PRV metrics as approximations of HR and HRV metrics obtained by wearable devices is by comparing them to a gold standard reference (Schuurmans et al., 2020). HR and HRV metrics obtained by an ECG device are often proposed as such a gold standard reference (e.g., van Lier et al., 2020).

## Figure 1

### *An Example of an ECG Signal and BVP Signal*



*Note.* On the left side of the vertical dotted line, the bio-electrical signal of the heart is shown, and on the right side, the systolic and diastolic pulse waves; PTT: pulse transit time.

### 1.1 Validity of E4 wristband-based PRV

One wearable device that its manufacturers termed fit for lab-based studies is the Empatica E4 wristband (i.e., E4 wristband; <https://www.empatica.com/en-eu/>). The E4 wristband has already been used to study the role of PR and/or PRV in a variety of lab-based research topics, such as stress (e.g., Mishra et al., 2020), music therapy (e.g., Rahman et al., 2021), and emotions (e.g., So et al., 2021).

To date, only a limited number of studies have examined the *validity* of the PR and PRV metrics obtained by the E4 wristband. For instance, Menghini et al. (2019) examined the validity of the E4 wristband's PR and PRV metrics as accurate approximations of ECG-based HR and HRV metrics under

various conditions (i.e., during seated position, paced breathing, orthostatic posture, slow walking, keyboard typing, Stroop test, speech preparation, public speech, and speech recovery). They showed that PR obtained by the E4 wristband and HR acquired with the gold standard ECG device was comparable across conditions. The PRV and HRV metrics in the time (SDNN and RMSSD) and frequency domain (LF and HF) were comparable during seated and paced breathing conditions, but these E4 wristband-based PRV metrics were less accurate approximations of the ECG-based HRV metrics during conditions that involved movement (i.e., slow walking) or cognitive performance (i.e., Stroop test). The authors demonstrated that this decrease in PRV accuracy in those conditions was likely due to hand/wrist movement (see Ryan et al., 2019, van Lier et al., 2020, for similar results). Furthermore, in contrast to the studies mentioned above, Schuurmans et al. (2020) observed that E4 wristband PRV-based RMSSD was not comparable to the ECG-based RMSSD during a seated-rest condition (see Ollander et al., 2016 for similar results). Additionally, McCarthy et al. (2016) simultaneously measured ECG using a portable ECG device and BVP using the E4 wristband during a 24h or 48h recording. They showed that in most cases (85%), signal quality assessed based on visual inspection (shape, stability, and noise) was comparable between the E4 wristband and the ECG device, although this was only the case during the less-active parts of the day (i.e., night and morning). Hence, previous studies assessing the validity of the PR and PRV metrics obtained with the E4 wristband seem to indicate that in seated or low activity conditions, the E4 wristband can provide comparable PR and PRV metrics as the HR and HRV metrics acquired with a gold standard ECG device, although for some indices (e.g., RMSSD) contradictory results have been obtained. However, inconsistencies in the statistical procedures used to assess the E4 wristband's validity (e.g., difference factor, visual inspection, and/or Bland-Altman plots; McCarthy et al., 2016; Ollander et al., 2016; Schuurmans et al., 2020), as well as the often small sample sizes (e.g.,  $N = 7$ ; Ollander et al., 2016), make it currently difficult to draw strong conclusions.

## **1.2 Validity of E4 wristband-based PRV obtained with ultra-short-term intervals**

As a gold standard, 5-minute intervals are stated as sufficient to measure PR/HR and PRV/HRV reliably (Task Force of the European Society of Cardiology and the North American Society of Pacing and Electrophysiology, 1996; from now on Task Force, 1996). However, lab experiments often consist of multiple short-interval events (e.g., "trials") occurring within one testing session (i.e., intervals < 5min; e.g., Choi et al., 2017; Lackner et al., 2013; Shen et al., 2017). A limited number of ECG studies have already compared the HR and HRV metrics obtained with UST intervals to those acquired with the gold standard 5min interval (e.g., Baek et al., 2015). For example, Shaffer et al. (2016) measured HR and HRV metrics via ECG with UST intervals (i.e., interval length = 10s, 20s, 30s, 60s, 90s, 120s, 180s, and 240s) in a seated-rest condition. They subsequently compared those HR and HRV metrics obtained with the UST intervals to those acquired with a 5min interval. They found that only HR across all UST intervals was highly correlated with the 5min HR. RMSSD and SDNN needed at least a 60s UST interval, and LF and HF power were only validly estimated with at least the 90s and 180s UST intervals, respectively (Shaffer & Ginsbeg, 2017; Shaffer et al., 2016; see also Baek et al., 2015). Munoz et al. (2015), on the other hand, showed that RMSSD and SDNN might already be validly estimated with a UST interval of 10s (see also Nussinovitch et al., 2011). However, to our knowledge, such a comparison between UST intervals and the gold standard 5min interval has never been studied for the E4 wristband's PR and PRV metrics. This is particularly important to establish the E4 wristbands' applicability as a PR and PRV recording device in lab studies comprising multiple fast trials.

Given this overview, the aim of the current study was twofold. First, we aimed to evaluate the validity of the E4 wristband's PR and PRV metrics as accurate approximations of HR and HRV by comparing it to a gold standard ECG device in a lab context. Second, we aimed to examine the validity

of the E4 wristband's UST interval measurement of PR and PRV to assess their usefulness for implementation in lab-based studies with short-interval events. We solely focused on a seated-rest condition in a lab context.

## 2. METHOD

### 2.1 Participants

A convenience sample of 79 undergraduates of the KU Leuven participated in this study in return for course credits. This sample was recruited in the context of a broader study on creative problem-solving (Stuyck et al., in preparation). In line with Quintana and Heathers (2014), several inclusion criteria (i.e., nonsmokers, body-mass index [BMI] <30, Beck depression inventory score <29, no cardiovascular or neurological medication use, no history of or current cardiopulmonary diseases, psychiatric disorders, and/or neurological disorders) and instructions concerning daily routines immediately preceding participation (i.e., no alcohol consumption the night before and day of the experiment, at least six hours sleep the night before the experiment, no caffeine consumption during the two hours before the experiment, no heavy meal consumption and no strenuous physical activity prior to the experiment) were used, as they might negatively influence the BVP and ECG recording.

We excluded seven participants based on technical issues with the ECG signal (i.e., incorrect electrode placement, disconnection of the electrodes during recording, or trigger interface unresponsiveness). In line with van Lier et al. (2020), we visually inspected the quality of the BVP signal morphology. We excluded participants if more than 20% of the recording interval contained an unstable BVP signal, indicating unusable data (see Appendix A for supplementary materials where supplementary material A depicts an example). This led to the exclusion of 26 participants. We also visually inspected the quality of the ECG signal morphology for abnormalities and stability (see supplementary material A for an example). This led to the additional exclusion of five participants (see Kumral et al., 2019; Shaffer & Ginsberg, 2017). After that, we identified outlying HRV observations by using the Tukey method (1997) of the median  $\pm$  three times the interquartile range (see Kumral et al., 2019 for a similar procedure). Based on this method, we excluded two participants with outlying PRV/HRV observations on two or more different PRV/HRV indices. Finally, we identified two participants with only outlying BVP-based and ECG-based LF. Therefore, we decided to exclude these participants from the statistical analyses involving LF solely<sup>1</sup>. This resulted in a final sample of 39 participants (mean age = 19,  $SD = 1.55$ , range = 17-24, 35 female). Our sample size was based on van Lier et al. (2020), who performed an a priori sample size calculation. Based on the large effect sizes observed in previous studies comparing the E4 wristband to a gold standard ECG device (i.e., minimum effect size of  $r = .72$ ; Schuurmans et al., 2020), and when comparing UST HRV to 5min recording intervals (i.e., minimum effect size of  $r = .76$ ; Munoz et al., 2015), our study with a sample of  $N = 39$  had an estimated power of .99 to detect such large effects (Campbell & Thompson, 2012). Before the start of the experiment, all participants provided written informed consent. The social and societal ethics committee of the KU Leuven approved this study (approval code G-2019 12 1929).

---

<sup>1</sup> We also performed the data analysis with the outliers included. The results obtained with this analysis are similar to the results of the main analysis described in the main text. This data analysis can be found in supplementary material B.

## 2.2 Assessment and Measurement

Before assessing the PRV and HRV indices crucial for this study, we note that participants completed six practice trials of a cognitive task, namely the compound remote associates test (CRAT; see <https://osf.io/snb3k/>). Once the PRV and HRV measurements relevant to this study were collected, participants performed this cognitive task as part of a broader study. As this task was not relevant to the current research questions and was only assessed *after* the critical PRV and HRV indices were collected, it will not be discussed further (for a detailed description of the CRAT instructions and experimental procedure, see <https://osf.io/4frcb/>).

## 2.3 Equipment

Participants were seated individually in a quiet, dimly lit room held at a constant temperature between 21° and 23°. They faced the computer monitor from approximately 60cm. A Dell Optiplex 3060 computer was used with a Dell 23.6-inch monitor.

### 2.3.1 Nexus-10 MKII ECG device

Nexus-10 MKII (Mind Media BV, Herten, the Netherlands) was used as the gold standard ECG recording device (CE-certified; 93/42/EC Annex XII). The device obtains the ECG signal in microvolts with a sampling rate of 256 Hz. Three pre-gelled Ag/AgCl electrodes were used. Following the modified Lead-II placement, these were attached to the upper body (see Kuipers et al., 2017). Namely, the negative electrode below the center of the right collarbone, the positive electrode on the lower left rib cage and the ground electrode below the left collarbone. Before placing the electrodes, the skin was cleaned with an alcohol pad.

### 2.3.2 Empatica E4 wristband

The Empatica E4 wristband (Empatica, Milan, Italy) is a certified medical wrist-worn device (CE-certified. No. 1876/MDD) that enables real-time multi-sensor data acquisition. The device allows the recording of four psychophysiological indices: BVP, acceleration, skin conductance, and skin temperature. The current study only used the BVP. The E4 wristband extracts BVP via a PPG sensor with two green and red photodiodes (LEDs). This BVP signal is acquired at a sampling rate of 64 Hz. Following the manufacturer's guidelines, participants wore the E4 wristband on their non-dominant hand to diminish the likelihood of motion artifacts.

## 2.4 Procedure

All participants were assessed between 9 am and 5 pm. Before entering the laboratory, the experimenter stressed that it was important to go to the toilet before testing if needed, as this might influence PRV/HRV data. After that, participants signed the informed consent. Their eligibility was assessed based on the inclusion criteria and adherence to the instructed daily routines. Participants were then given instructions on how to attach the ECG device's electrodes, which was accompanied by an image depicting the exact modified Lead-II placement. They also attached the E4 wristband to the wrist of their non-dominant hand. The experimenter visually checked if both devices were attached correctly. Participants then took place at the test computer. The experimenter explained that it was important to remain in an upright seated position without crossing their legs and to minimize sudden movements as much as possible (e.g., arm stretching, excessively coughing). It was explained that deviating from these

instructions might cause BVP and ECG data acquisition issues. Participants then completed the six practice trials of the CRAT computer task. The experimenter explained that there would be a 10min resting period after these practice trials. During this 10min resting period, participants were instructed to relax and control their emotions to minimize the likelihood of engaging in ruminative or emotionally valenced thoughts that might affect the PRV/HRV recording. To minimize recording-awareness-related changes in their behavior, we told participants that we were mainly interested in their (later) performance during the cognitive task. This 10min period consisted of a 5min acclimatization period and a subsequent 5min baseline period. It is this 5min baseline period that was used in this study to validate the E4 wristband. After the 10min resting period, participants completed the CRAT.

## 2.5 Data preprocessing

The E4 wristband and ECG device data were synchronized by relying on the computer's clock time (see Menghini et al., 2019; Milstein & Gordon, 2020). The Empatica software responsible for the E4 wristband recording synchronizes to the computer's clock time. Therefore, we created a timestamp of the computer's clock time while at the same time sending a trigger to the ECG device at the onset of the experiment. Subsequently, we matched the clock time provided by this timestamp, represented by the trigger in the ECG device, to that of the clock time registered by the Empatica software, thereby creating a synchronous time point of reference for the ECG device (via the trigger) and the E4 wristband. To ensure tight time synchronization between the two devices, we visually double-checked the concordance between the PPI and RRI time series.

Both the data obtained by the ECG device (i.e., ECG signal) and the E4 wristband (i.e., BVP signal) were preprocessed using Kubios premium (v. 3.4.2; Tarvainen et al., 2014 and see Tarvainen et al., 2020). The IBIs were calculated by determining the time interval between two r-peaks for the ECG signal and between two fiducial points on the systolic pulse wave for the BVP signal. To correct potential artifacts in the IBI time series, we used the automatic artifact correction algorithm of Kubios (see Lipponen & Tarvainen, 2019). All detected artifacts are subsequently replaced with IBIs based on cubic spline interpolation. To accommodate the non-stationarity of the IBI time series, Kubios deploys a detrending procedure by defining an a priori smoothing parameter (cut-off frequency 0.035 Hz; see Tarvainen et al., 2002). Additionally, all ECG and BVP signals were visually inspected for abnormal signals, unstable recording epochs in the recording interval, missed p-waves/r-peaks and missed artifacts by the algorithms that might influence the PRV/HRV data (e.g., supraventricular extrasystole; see Kumral et al., 2019 for a similar procedure). In case of any ECG or BVP signal abnormalities missed by the algorithms, we applied a manual correction to the signal (e.g., marking an unstable signal epoch as a to-be-excluded noise segment or adding missed p-waves/r-peaks; see supplementary material C for an overview of the percentage of corrected IBIs and noise-free ECG/BVP recording intervals).

### 2.5.1 Validity of E4 wristband-based PRV

Besides mean PR and HR (i.e., expressed in beats per minute; bpm), we used several PRV and HRV metrics to compare both recording devices. From the time domain, RMSSD and SDNN (both expressed in ms) were used and, from the frequency domain, the absolute power in the LF and HF bands (both expressed in  $\text{ms}^2$ ) were used (similar to Menghini et al., 2019; Schuurmans et al., 2020; van Lier et al., 2020). We did not include the LF/HF ratio, as it remains an ambiguous HRV metric (Heathers, 2014).



### 2.5.2 Validity of E4 wristband-based PRV obtained with UST intervals

To compare the E4 wristband's PR and PRV obtained with the UST intervals to those same metrics acquired with the 5min interval, we used RMSSD, SDNN, and HF (similar to Baek et al., 2015; Munoz et al., 2015; Salahuddin et al., 2007). The estimation of power in the HF band requires at least 10 oscillations (i.e., 70s; Task Force, 1996). As BVP-based PRV estimation is considered less stable than ECG-based HRV estimation, we argue that it is unlikely that this bare minimum of 70s will suffice for an accurate estimation of HF. Therefore, we chose only to use a UST interval of 120s to estimate HF. We did not assess LF as it requires at least 10 oscillations (i.e., 250s), making it impossible to estimate it with UST HRV recordings (Pecchia et al., 2018; Task Force, 1996). All six UST intervals were segmented from the 5min interval. We randomly extracted three nonoverlapping 10s segments from the 5min recording for each participant. Note that the extraction was nonsequential, in the sense that T1 did not necessarily precede T2, etc. These 10s segments were used independently to assess the PR and PRV metrics and calculate an average value of the PR and PRV metric across those three 10s segments. As the 10s intervals only contain very few IBIs, we only accepted a 100% stable BVP signal without corrected IBIs. If this was not the case, we rejected the 10s interval and randomly selected a new one. This process we repeated until we obtained three clean, valid 10s intervals. The UST interval of 30s existed of the first 30s of the 5min interval, and the 120s UST interval existed of the 120s interval after those first 30s. This ensures their independence in validly estimating PR and PRV (see Munoz et al., 2015 for a similar procedure).

Our first UST analysis assessed the E4 wristband's internal consistency concerning PR and PRV measured at UST intervals. This is vital as it shows that measuring PR and/or PRV at these time scales can be surrogates for their 5min estimation with the same device, thereby taking variability specific to BVP into account. However, the ECG 5min recording can still be considered as the best approximation of the genuine mean HR and HRV. Therefore, in a second UST analysis, we compared the PR and PRV estimation of the UST intervals to the 5min ECG-based estimation of HR and HRV. In the results section, we summarized the results from this second UST analysis, and full details can be found in supplementary material D.

## 2.6 Statistical analysis

A *three-step hierarchical procedure* was used to assess the validity of PRV metrics obtained with the E4 wristband and UST recording intervals (Pecchia et al., 2018, Shaffer et al., 2020, van Lier et al., 2020).

### 2.6.1 Step 1, PRV metric selection

A Pearson product-moment correlation coefficient was used to examine the association strength between the proxies (PRV metrics obtained with the E4 wristband or UST interval) and the gold standard measurement (HRV metrics acquired with the ECG device or PRV metrics acquired with the 5min interval). We followed the procedure of Menghini et al. (2019) and the recommendations of Pecchia et al. (2018). We used a correlation coefficient cut-off of  $r = .70$  to identify the viable proxies of the gold standard measurement. Only the proxies that showed a correlation with the gold standard greater than or equal to  $r = .70$  were retained for further analysis.

### 2.6.2 Step 2, PRV metric validity

Subsequently, a Bland-Altman plot was created for the viable proxies selected in step 1. Here, the differences between the proxy and the gold standard measurement are plotted against the gold standard measurement (see Giavarina, 2015; Menghini et al., 2021, for an in-depth explanation). We plotted against this gold standard instead of against the mean of the proxy and the gold standard, as suggested by Bland and Altman (1995), because in the current study, the gold standard is expected to be the best approximation of the genuine PRV and HRV. Therefore, it is also expected to have lower error variance and bias and, thus, to be better suited than the mean of the proxy and gold standard as a reference to plot against (see Krouwer, 2008; Munoz et al., 2015, for similar argumentation).

First, the mean of this difference was analyzed. This mean reflects the bias in measurement, as a perfect agreement between the proxy and the gold standard measurement would be reflected by a mean difference of zero. This bias can reflect a tendency of E4 wristband or UST interval to over/underestimate the proxy relative to the gold standard. The 95% confidence intervals of the bias (i.e., 95% CI) were used to assess this (Menghini et al., 2021). If zero is below/above the lower/upper bounds of this 95% CI, there is an indication of a tendency to over/underestimate the proxy relative to the gold standard.

Second, it was measured to what extent the observed differences between the proxy and the gold standard were within acceptable limits of agreement. The limits of agreement (LOA) are represented by the mean of the differences (bias)  $\pm 1.96SD$  (Menghini et al., 2021). The bounds of these LOA mark the inclusion of 95% of the observed differences (i.e., 95% LOA). To measure if those 95% LOA can be considered acceptable limits in which the majority of the differences between the estimated proxy and gold standard lie, *a priori* LOA need to be determined before constructing the Bland-Altman plot (Giavarina, 2015). If the LOA are within the *a priori* LOA, the bounds marking the inclusion of 95% of the observed differences are within the limits of a maximum acceptable deviation of the gold standard. This would imply that the proxy sufficiently agrees with the gold standard.

The determination of the *a priori* LOA requires the consideration of several features of the research (Giavarina, 2015; van Lier et al., 2020). For instance, the clinical necessity (e.g., the aimed-for diagnostic accuracy of the proxy), biological features of the sample being studied, the aim of the study, the time interval used to estimate PR and HR, and PRV and HRV (e.g., seconds, minutes or hours), and the parameters being studied (e.g., HR and RMSSD; van Lier et al., 2020). In the literature, mainly two *a priori* LOA have been proposed: an *a priori* LOA of 150% (e.g., Menghini et al., 2019) and an *a priori* LOA of 110% (e.g., van Lier et al., 2020). For example, an average 50ms RMSSD measured with the gold standard leads to an *a priori* 150% LOA of  $\pm 25\text{ms}$  (i.e., *lower bound* =  $\frac{\bar{x} \text{ gold standard} \times 150}{100} - \bar{x} \text{ gold standard}$ ; *upper bound* =  $\frac{\bar{x} \text{ gold standard} \times 150}{100}$ ). Thus, the 95% LOA should lie within these limits of the maximum acceptable deviation of the gold standard of  $\pm 25\text{ms}$ . Although an *a priori* LOA of 110% is considerably stricter than one of 150%, this stricter *a priori* LOA is generally used to assess the utility of medical equipment (Advancement of Medical Instrumentation, 2002; van Lier et al., 2020). In this case, such a strict *a priori* LOA seems justified as medical diagnoses should depend on highly accurate derived metrics. However, as the current study aimed to validate the proxies obtained with E4 wristband and UST intervals within a lab research context with a maximal recording length of 5min, we argue that the margin of error between the proxy and gold standard can be less strict. Therefore, we used an *a priori* LOA of 150% calculated as in the example above, consistent with similar previous studies (Charlot et al., 2009; Menghini et al., 2019, Pichon et al., 2006).

Several assumptions need to be considered for the construction of the Bland-Altman plot. The bias and the 95% LOA in the Bland-Altman plot are crucial for its interpretation. However, their accurate representation also depends on the absence of a positive/negative association between the

observed differences and the gold standard value (i.e., proportional bias), an equal spread of the observed differences for each gold standard value (i.e., homoscedasticity), and the normal distribution of the observed differences between the proxy and the gold standard (i.e., normality). If one of these assumptions is not met, the above-explained representations of the bias and 95% LOA no longer hold. In that case, the bias and 95% LOA need to be represented via other calculations taking into account the violations (see supplementary material E for the precise calculations). To identify and accommodate any violations of the assumptions in our data, we followed the procedure as described by Altman and Bland (1983), Bland and Altman (1999), Euser et al. (2008), and Menghini et al. (2021). Regardless of the type of calculation used to represent the bias and its 95% LOA, we always represented the observed differences, the gold standard values, the bias, the 95% LOA, and the a priori 150% LOA on their original scale to enhance interpretability.

### **2.6.3 Step 3, magnitude of difference**

In the final step of the procedure, we calculated Cliff's delta ( $d$ ) effect size to estimate the degree of overlap between the distributions of the proxy and the gold standard. This nonparametric effect size is less vulnerable to skewed, non-normal, and heteroscedastic data than, for example, Cohen's  $d$  (Romano, 2006). Its absolute values range from 0 to 1. We interpreted Cliff's  $d$  effect size as follows:  $d \leq .15$  = negligible,  $.15 < d \leq .33$  = small,  $.33 < d \leq .5$  = medium, and  $d > .5$  = large (see Romano, 2006).

We used the open-source *R* language and environment to perform statistical analysis (R Core Team, 2021). For the Bland-Altman plots, we used the source code provided by Menghini et al. (2021). We adjusted this source code to the specifics of the current study. To compute Cliff's  $d$  and its corresponding 95% confidence intervals, we used the "effsize" package (Torchiano, 2020). All R code can be found on the Open Science Framework (<https://osf.io/4frcb/>).

## **3. RESULTS**

### **3.1 Validity of E4 wristband-based PRV**

#### **3.1.1 Step 1, PRV metric selection**

Table 1 depicts the correlation coefficients and the descriptive statistics of the mean PR and HR, the time (RMSSD and SDNN), and frequency (LF and HF) domain PRV and HRV metrics. Correlations between both devices surpassing the cut-off of  $r = .70$  were found for HR, RMSSD, SDNN, LF, and HF. Therefore, the mean PR and all PRV metrics were included in step 2 of the three-step hierarchical procedure.

**Table 1.**

*Correlation Coefficients and the Descriptive Statistics of PR/HR and PRV/HRV Metrics obtained with E4 and ECG*

Metrics	$r$ [95% CI]	$M(SD)$ E4	$M(SD)$ ECG
PR/HR	.9989[.9982, .9997]	82.06(7.13)	82.20(7.10)
RMSSD	.92[.88, .97]	39.33(10.65)	33.34(11.43)
SDNN	.98[.97, 1.00]	44.50(13.29)	42.92(13.91)
LF	.96[.94, .99]	1030.22(674.02)	1066.43(704.99)
HF	.95[.91, 1.00]	764.77(518.23)	637.73(450.89)

*Note.* PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); LF, low frequency (ms<sup>2</sup>); HF, high frequency (ms<sup>2</sup>);  $r$ , Pearson correlation coefficient; 95% CI, 95% confidence intervals.

### 3.1.2 Step 2, PRV metric validity

Table 2 presents the results, and Figure 2 presents the Bland-Altman plots (see supplementary material E for the assumptions handling). Concerning mean PR/HR, there was a slight tendency of the E4 wristband to underestimate it relative to its estimation by the ECG device. Namely, the upper bound of the bias's 95% CI (i.e., -0.006) was just below zero. However, the 95% LOA was within the *a priori* 150% LOA, with 100% of the observed differences lying within the *a priori* 150% LOA bounds. These findings illustrated that the deviations between the mean PR obtained with the E4 wristband and the mean HR acquired with the ECG device were within the maximum acceptable deviation limits.

For RMSSD, the 95% CI of its bias was entirely above zero, indicating that the E4 wristband tended to overestimate RMSSD relative to the same metric acquired with the ECG device. For SDNN, on the other hand, proportional bias was detected (i.e., a negative association between the observed differences and the ECG-device's estimation of SDNN; see Figure 2). Specifically, SDNN was overestimated by the E4 wristband for observations of ECG-based SDNN lower than 50ms, whereas the bias was non-significant for higher ECG-based SDNN observations. For RMSSD and SDNN, the 95% LOA was within the *a priori* 150% LOA, with 97% and 100% of the observed differences lying within the *a priori* 150% LOA bounds, respectively. Therefore, the E4 wristband's estimation of RMSSD and SDNN may be considered sufficiently in agreement with the same metrics estimated by the ECG device.

For LF and HF, a tendency to overestimate these metrics by the E4 wristband relative to the ECG device was observed. Namely, in both cases, the entire 95% CI of the bias was above zero. For LF, the 95% LOA were within the *a priori* 150% LOA for 89% of the ECG-based LF observations, with acceptable 95% LOA for observations of ECG-based LF lower than 1979ms<sup>2</sup>. Here, 97% of the observed differences were lying within the *a priori* 150% LOA bounds. As such, the E4 wristband estimation of LF was in sufficient agreement with its estimation by the ECG device. However, for HF, the 95% LOA were within the *a priori* 150% LOA for 46% of the ECG-based HF observations, with acceptable 95% LOA for observations of ECG-based HF lower than 432ms<sup>2</sup>. Here, only 87% of the observed differences were within the *a priori* 150% LOA. This finding illustrated that the E4 wristband's estimation of HF deviated too much from its gold standard estimation by the ECG device.

Consequently, the mean PR, RMSSD, SDNN, and LF were retained for the last step of the three-step hierarchical procedure.

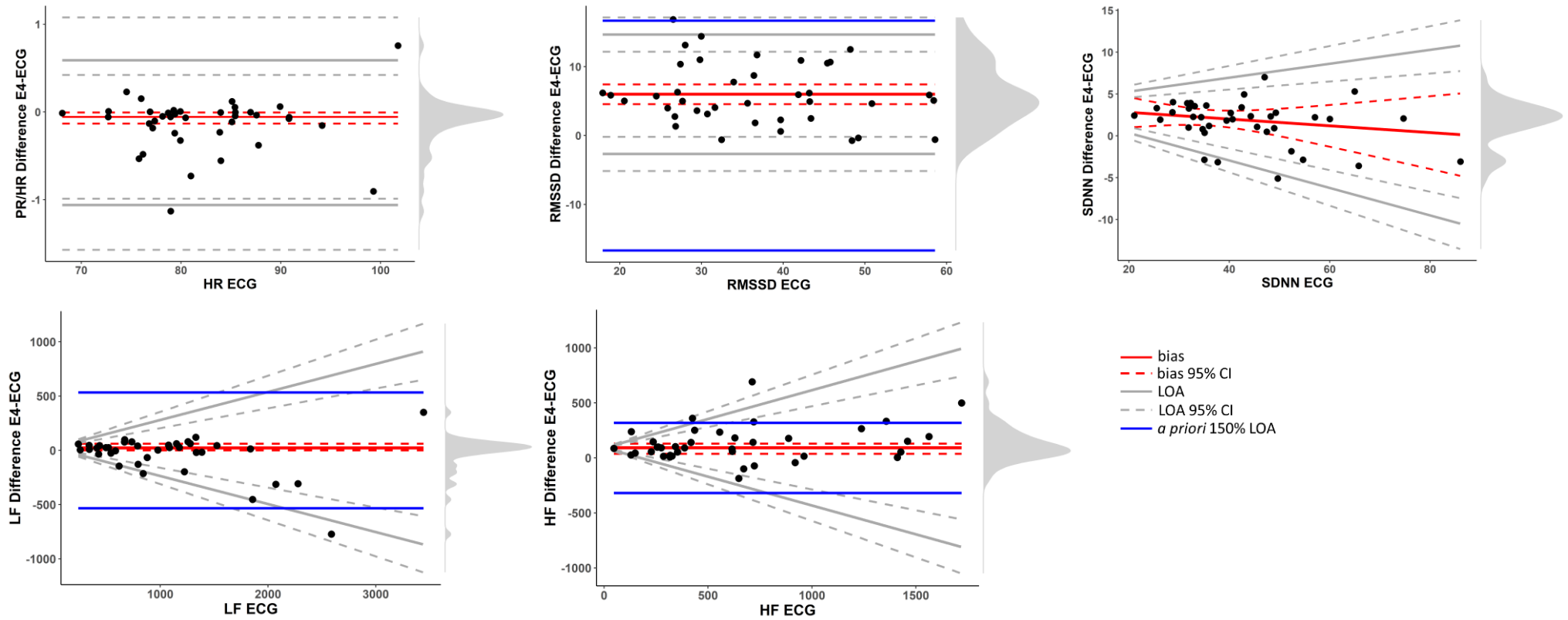
**Table 2.***Bland-Altman Analysis: Bias, 95% LOA, and A Priori 150% LOA*

Metrics	Bias [95% CI]	Lower LOA [95% CI]	Upper LOA [95% CI]	<i>A priori</i> 150% LOA
PR/HR*	-0.06 [-0.13, -0.006]	-1.06 [-1.57, -0.99]	0.59 [0.42, 1.08]	-41.10, 41.10
RMSSD*	5.99 [4.56, 7.42]	-2.66 [-5.15, -0.19]	14.65 [12.17, 17.13]	-16.67, 16.67
SDNN*	3.62 - 0.04*GS [2.10, 6.49], [-0.11, -0.004]	Bias - 0.12*GS [0.09, 0.16]	Bias + 0.12*GS [0.09, 0.16]	-21.46, 21.46
LF*	21.69 [0.13, 61.32]	Bias - 0.26*GS [0.18, 0.33]	Bias + 0.26*GS [0.18, 0.33]	-533.21, 533.21
HF	91.60 [37.00, 129.30]	Bias - 0.52*GS [0.38, 0.66]	Bias + 0.52*GS [0.38, 0.66]	<b>-318.86, 318.86</b>

*Note.* PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); LF, low frequency (ms<sup>2</sup>), HF, high frequency (ms<sup>2</sup>); *PR/HR distribution characteristics*, homoscedastic and non-normal which a log transform could not alleviate; *RMSSD distribution characteristics*, homoscedastic and normal; *SDNN distribution characteristics*, heteroscedastic and non-normal which a log transform could not alleviate; *LF distribution characteristics*, heteroscedastic and non-normal which log transformation could not alleviate; *HF distribution characteristics*, heteroscedastic and non-normal which a log transform could not alleviate (see supplementary material E for the assumption handling); Median; 2.5 percentile; 97.5 percentile;  $\beta_0$ , the intercept;  $\beta_1$ , the slope coefficient; antilog, the antilog slope value; GS, the gold standard value; LOA, limits of agreement; 95% CI, 95% confidence intervals; *a priori* 150% LOA, *a priori* 150% limits of agreement; bold typeface indicates when the 95% LOA is outside the *a priori* 150% LOA; \* indicates the proxies that were retained for the following step in the three-step hierarchical procedure.

Figure 2

Bland-Altman Plots for Mean PR/HR and PRV/HRV Metrics



Note. PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); LF, low frequency ( $\text{ms}^2$ ), HF, high frequency ( $\text{ms}^2$ ); density distribution of the differences at the right side of each plot; 95% CI, 95% confidence interval; LOA, limits of agreement; *a priori* LOA of 150% (only displayed if close to 95% LOA).

### ***3.1.3 Step 3, the magnitude of difference***

The observed effects sizes were negligible for mean PR/HR (Cliff's  $d = -.03$ , 95% CI [-.28, .23]). For the time-domain metrics, the effect sizes were small for RMSSD (Cliff's  $d = .30$ , 95% CI [.04, .52]) and negligible for SDNN (Cliff's  $d = .10$ , 95% CI [-.15, .35]). For the frequency-domain metric LF, the effect size was negligible (Cliff's  $d = -.01$ , 95% CI [-.27, .25]). These findings show that the E4 wristband provided comparable estimates of mean PR, RMSSD, SDNN, and LF to their same mean HR and HRV estimation with the ECG device.

## **3.2 Validity of E4 wristband-based PRV obtained with UST intervals**

### ***3.2.1 E4 wristband UST interval PRV versus 5min E4 wristband-based PRV***

**3.2.1.1 Step 1, PRV metric selection.** The correlation coefficients and descriptive statistics of the mean PR and the PRV metrics of the UST intervals are presented in Table 3. Correlations between UST intervals and 5min intervals surpassing the cut-off of  $r = .70$  were found for all UST intervals estimating PR, RMSSD, SDNN, and HF, except for two of the three UST intervals of 10s for RMSSD and SDNN. Hence, all UST intervals estimating PR and PRV surpassing the cut-off were included in the second step of the three-step hierarchical procedure.

**Table 3.**

*Correlation Coefficients and Descriptive Statistics of the PR and PRV Metrics for the UST and 5min Recordings Obtained with the E4 Wristband*

Metrics	<i>r</i> (95% CI)	<i>M</i> ( <i>SD</i> )
<i>PR</i>		
5min		82.06(7.13)
10s T1	.90[.83, .98]	82.22(7.80)
10s T2	.89[.82, .98]	82.00(8.00)
10s T3	.93[.88, .98]	81.27(8.51)
Average of 10s	.96[.93, .99]	81.83(7.67)
30s	.90[.82, .99]	82.19(7.40)
120s	.98[.96, 1.00]	82.05(6.74)
<i>RMSSD</i>		
5min		39.33(10.65)
10s T1	.44[.15, .74]	39.64(15.55)
10s T2	.70[.53, .85]	36.78(14.95)
10s T3	.67[.51, .82]	36.27(14.42)
Average of 10s	.80[.64, .96]	37.56(11.26)
30s	.79[.67, .90]	37.72(12.23)
120s	.93[.90, .97]	39.80(10.94)
<i>SDNN</i>		
5min		44.50(13.29)
10s T1	.47[.20, .77]	49.24(20.85)
10s T2	.61[.42, .80]	38.84(18.22)
10s T3	.76[.57, .98]	39.93(17.59)
Average of 10s	.79[.63, .99]	42.67(14.47)
30s	.78[.65, .94]	42.83(15.20)
120s	.93[.88, .99]	45.59(14.44)
<i>HF</i>		
5min		764.77(518.23)
120s	.84[.75, .95]	725.48(511.31)

*Note.* PR, mean pulse rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); T, time interval; *r*, Pearson correlation coefficient; 95% CI = 95% confidence interval.

**3.2.1.2 Step 2, PRV metric validity.** The results of the Bland-Altman analysis are depicted in Table 4, and exemplary Bland-Altman plots are shown in Figure 3 (see supplementary material E for the assumptions handling). The bias included zero for all UST intervals estimating mean PR in its 95% CI, except for the UST interval of 120s. Here, proportional bias was detected (i.e., a negative association between the observed differences and the 5min interval estimation of mean PR; see Figure 3). Specifically, PR was slightly underestimated by the UST interval of 120s for observations of 5min-based PR higher than 96 bpm, whereas the bias was non-significant for lower 5min-based PR values. For all other UST intervals, no tendency to over/underestimate mean PR was found relative to its same estimation with the 5min interval. For all UST intervals' estimation of mean PR, the 95% LOA was



inside the *a priori* 150% LOA, thereby showing that the UST intervals provided comparable mean PR values to the same value acquired with the 5min interval. In all cases, 100% of the data lay within the *a priori* 150% LOA.

For all the UST intervals estimating RMSSD, the bias included zero in its 95% CI. Therefore, for all these UST intervals, no tendency to over/underestimate RMSSD relative to the same estimation with the 5min interval was found. For the UST interval of 30s, the 95% LOA was within the *a priori* 150% LOA for 85% of the 5min-based RMSSD observations, with acceptable 95% LOA for observations of 5min-based RMSSD lower than 51ms. Here, 97% of the observed differences were lying within the *a priori* 150% LOA bounds. For the UST intervals of 120s, the 95% LOA were inside the *a priori* 150%, with 100% of the observed differences inside the *a priori* 150% LOA boundaries. For the average of the three UST intervals of 10s, the lower bound of the 95% LOA was borderline outside the lower bound of the *a priori* 150% LOA (see Figure 3), with 97% of the observed differences within the *a priori* 150% LOA boundaries. Therefore, the RMSSD estimated with the average of the three UST intervals of 10s, the UST intervals of 30s and 120s were in sufficient agreement with its estimation by the 5min interval. However, for the UST intervals of 10s (i.e., T2), the 95% LOA was outside the *a priori* 150% LOA, with only 90% of the observed differences within the *a priori* 150% LOA boundaries. Therefore, the estimation of RMSSD with this UST interval of 10s deviated too much from the same value acquired with the 5min interval, indicating insufficient agreement.

For all UST intervals estimating SDNN, the bias included zero in its 95% CI, except for the UST intervals of 10s. Here, the upper bound of the 95% CI was below zero, indicating a tendency of this UST interval to underestimate SDNN relative to its estimation with the 5min interval. For all other UST intervals estimating SDNN, no tendency to over/underestimate it was found relative to its same estimation with the 5min interval. The 95% LOA was within the *a priori* 150% LOA for the average of the three UST intervals of 10s, the UST interval of 30s, and 120s. For these UST intervals, 97%, 97%, and 100% of the observed differences were inside the *a priori* 150% LOA boundaries, respectively. This result showed that the SDNN estimated by the average of the three UST intervals of 10s, the UST interval of 30s and 120s was in sufficient agreement with its same estimation by the 5min interval. For the UST interval of 10s, the 95% LOA were within the *a priori* 150% LOA for 8% of the 5min-based SDNN observations, with acceptable 95% LOA for observations of 5min-based SDNN lower than 31ms. Here, only 90% of the observed differences were within the *a priori* 150% LOA boundaries. As such, the SDNN estimated with the UST interval of 10s deviated too much from its estimation with the 5min interval, thereby showing insufficient agreement.

For the UST interval of 120s estimating HF, the upper bound of the bias's 95% CI was below zero, indicating a tendency of this UST interval to underestimate HF relative to its same estimation with the 5min interval (see Figure 3). Furthermore, the 95% LOA were within the *a priori* 150% LOA for 8% of the 5min-based HF observations, with acceptable 95% LOA for observations of 5min-based HF lower than 274ms<sup>2</sup>. Here, only 82% of the observed differences were within the *a priori* 150% LOA boundaries. Therefore, the HF estimated with the UST interval of 120s deviated too much from its estimation with the 5min interval, indicating insufficient agreement.

Consequently, all UST intervals estimating mean PR, the average of the three UST intervals of 10s, the UST interval of 30s and 120s estimating RMSSD and SDNN, were further analyzed in the final step of the three-step hierarchical procedure.

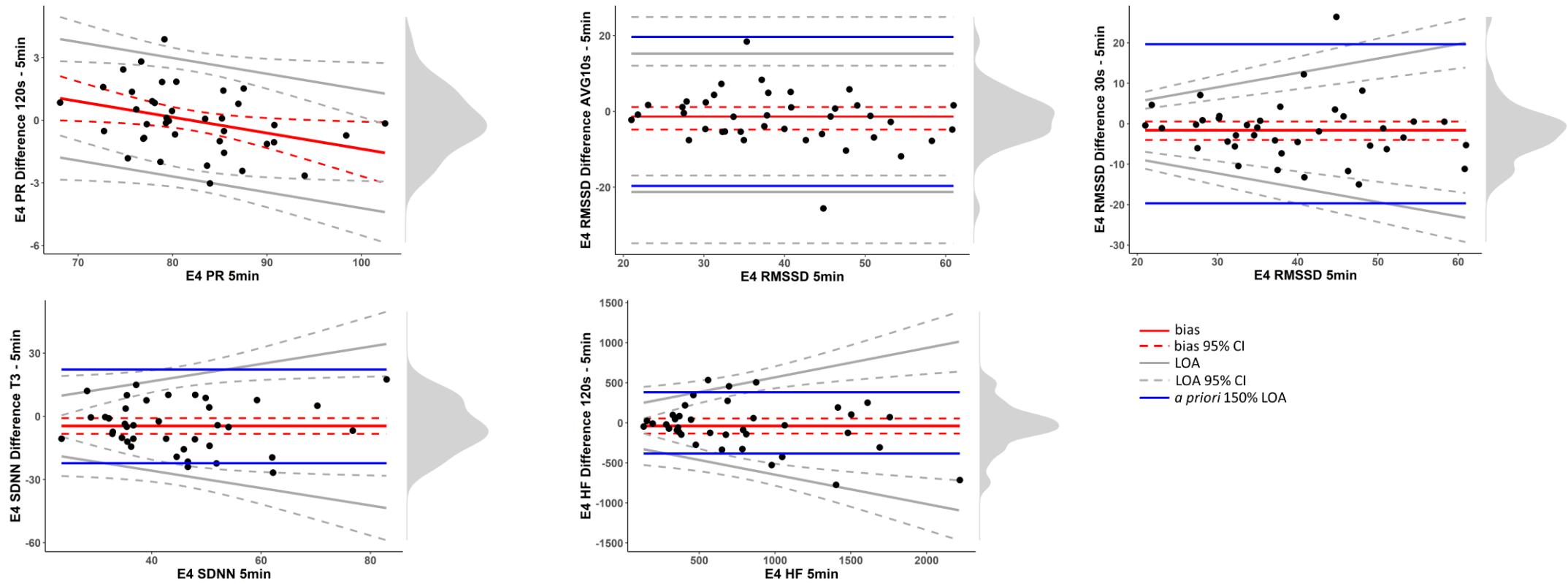
**Table 4.***UST Recording's Bland-Altman Analysis: Bias and 95% LOA*

Metrics	Bias [95% CI]	Lower LOA [95% CI]	Upper LOA [95% CI]	<i>A priori</i> 150% LOA
<i>PR</i>				
10s T1*	0.17 [-0.93, 1.25]	-6.42 [-8.31, -4.54]	6.75 [4.86, 8.64]	-41.03, 41.03
10s T2*	-0.06 [-1.22, 1.10]	-7.08 [-9.09, -5.07]	6.96 [4.95, 8.97]	-41.03, 41.03
10s T3*	-0.79 [-1.84, 0.26]	-7.15 [-8.98, -5.33]	5.57 [3.75, 7.39]	-41.03, 41.03
Average of 10s*	-0.23 [-0.93, 0.47]	-4.47 [-5.69, -3.26]	4.02 [2.80, 5.23]	-41.03, 41.03
30s*	0.13 [-0.91, 1.18]	-6.19 [-8.00, -4.38]	6.46 [4.65, 8.27]	-41.03, 41.03
120s*	<b>6.22 - 0.08*GS</b> [0.64, 11.80], [-0.14, -0.01]	Bias - 1.96*1.45 [1.18, 1.79]	Bias + 1.96*1.45 [1.18, 1.79]	-41.03, 41.03
<i>RMSSD</i>				
10s T2	<b>-4.58</b> [-6.77, 0.82]	<b>-25.15</b> [-38.02, -22.10]	<b>23.96</b> [23.75, 34.21]	<b>-19.67, 19.67</b>
Average of 10s*	<b>-1.33</b> [-4.79, 1.15]	<b>-21.24</b> [-34.77, -16.89]	<b>15.26</b> [12.06, 24.91]	<b>-19.67, 19.67</b>
30s*	-1.62 [-4.02, 0.56]	Bias - 0.36*GS [0.26, 0.45]	Bias + 0.36*GS [0.26, 0.45]	-19.67, 19.67
120s*	0.47 [-0.85, 1.79]	-7.49 [-9.77, -5.21]	8.43 [6.15, 10.71]	-19.67, 19.67
<i>SDNN</i>				
10s T3	-4.57 [-8.31, -0.84]	Bias - 2.46*(1.92 + 0.17*GS) [-5.20, 9.03], [0.01, 0.32]	Bias + 2.46*(1.92 + 0.17*GS) [-5.20, 9.03], [0.01, 0.32]	<b>-22.25, 22.25</b>
Average of 10s*	-1.83 [-4.76, 1.10]	-19.55 [-24.62, -14.47]	15.88 [10.81, 20.96]	-22.25, 22.25
30s*	-1.67 [-4.77, 1.42]	-20.38 [-25.74, -15.02]	17.04 [11.68, 22.39]	-22.25, 22.25
120s*	1.09 [-0.66, 2.84]	-9.50 [-12.54, -6.47]	11.68 [8.65, 14.71]	-22.25, 22.25
<i>HF</i>				
120s	-39.29 [-132.9, -54.30]	Bias - 2.46*(98.10 + 0.15*GS) [7.07, 189.10], [0.05, 0.25]	Bias + 2.46*(98.10 + 0.15*GS) [7.07, 189.10], [0.05, 0.25]	<b>-382.39, 382.39</b>

*Note.* PR, mean pulse rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); T = time interval; *PR UST intervals distribution characteristics*, all were homoscedastic and normally distributed with only the UST interval of 120s showing proportional bias; *RMSSD UST distribution characteristics*, all were homoscedastic with T2, the average of the three UST interval of 10s, and the UST interval of 30s displaying non-normality for which log transformation only alleviated non-normality for the UST interval of 30s; *SDNN UST intervals distribution characteristics*, all were homoscedastic and normally distributed except for the UST interval of 10s which was heteroscedastic; *HF UST interval distribution characteristics*, data was heteroscedastic and normally distributed (see supplementary material E for the assumption handling); **Median; 2.5 percentile; 97.5 percentile;  $\beta_0$ , the intercept;  $\beta_1$ , the slope coefficient; Bias SD, the SD of the residuals of the proportional bias model; antilog, the antilog slope value**; GS, the gold standard value; LOA = 95% limits of agreement; 95% CI = 95% confidence interval; *A priori* 150% LOA, *a priori* 150% limits of agreement; Bold typeface = the 95% LOA is outside the *a priori* 150% LOA; \* indicates the proxies that were retained for the following step in the three-step hierarchical procedure.

Figure 3

Examples of Bland-Altman Plots Comparing Mean PR and PRV Metrics Obtained with UST Intervals by the E4 Wristband Versus Their 5min Recording with the E4 Wristband



Note. PR, mean pulse rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); AVG10s, average of the three 10s intervals; T3, third UST interval of 10s; density distribution of the differences at right side plot; 95% LOA, 95% limit of agreement; 95% CI, 95% confidence intervals; the *a priori* LOA of 150% only displayed if close to 95% LOA.

**3.2.1.3 Step 3, magnitude of difference.** Regarding mean PR, the effect sizes were negligible for all the UST intervals: T1 (Cliff's  $d = -.21$ , 95% CI [-.27, .23]), T2 (Cliff's  $d = .03$ , 95% CI [-.22, .28]), T3 (Cliff's  $d = .11$ , 95% CI [-.15, .35]), the average of three 10s intervals (Cliff's  $d = .03$ , 95% CI [-.22, .27]), 30s interval (Cliff's  $d = -.01$ , 95% CI [-.26, .24]) and 120s interval (Cliff's  $d = -.01$ , 95% CI [-.27, .24]. Concerning RMSSD, negligible effect sizes were found for the average of the three UST intervals of 10s (Cliff's  $d = .10$ , 95% CI [-.15, .34]), the UST interval of 30s (Cliff's  $d = .12$ , 95% CI [-.13, .36], and the UST interval of 120s (Cliff's  $d = -.04$ , 95% CI [-.29, .22]). Lastly, for the UST intervals estimating SDNN negligible effect sizes were found for the average of the three UST intervals of 10s (Cliff's  $d = .10$ , 95% CI [-.16, .34]), the UST interval of 30s (Cliff's  $d = .04$ , 95% CI [-.22, .29]), and the UST interval of 120s (Cliff's  $d = -.04$ , 95% CI [-.29, .22]). These results show that the UST intervals specified above provide estimates of mean PR, RMSSD and SDNN comparable to the same value estimated with a 5min interval.

### 3.2.2 E4 wristband UST interval PRV versus 5min ECG-based HRV

**3.2.2.1 Step 1, PRV metric selection.** For all UST intervals, strong correlations ( $r > .70$ ) were observed between their estimation of PR and PRV and the 5min ECG-based estimation of HR and HRV, except for the three UST intervals of 10s for RMSSD and two of the three UST intervals of 10s for SDNN.

**3.2.2.2 Step 2, PRV metric validity.** Considering only those UST intervals that survived step 1, all UST intervals' estimations of PR were in sufficient agreement with the 5min ECG-based estimation of HR. Concerning RMSSD, only its estimation with the UST interval of 30s was in sufficient agreement with the same estimation with the 5min ECG recording. For the average of the three UST intervals of 10s and the UST interval of 120s, only 92% of the observed differences were within the *a priori* 150% LOA. Regarding SDNN, only its estimation with the average of the three UST intervals of 10s and the UST interval of 30s and 120s were in sufficient agreement with their same estimation with the 5min ECG recording. For the SDNN estimation with the UST interval of 10s, only 90% of the observed differences were within the *a priori* 150% LOA. Concerning HF, its estimation with the UST interval of 120s agreed insufficiently with the same estimation with the 5min ECG recording.

**3.2.2.3 Step 3, magnitude of difference.** Considering only those UST intervals that survived step 2, negligible effect sizes were found between the UST intervals estimation of PR and SDNN and the 5min ECG-based estimation of HR and SDNN. A small effect size was found for the UST interval of 30s estimating RMSSD (see supplementary material D for the detailed statistical results).

## 4. DISCUSSION

The current study aimed to validate the mean PR and PRV metrics obtained with the E4 wristband as approximations of mean HR and HRV metrics by comparing them to the mean HR and HRV metrics acquired with a gold standard ECG device. Moreover, we assessed the time scales at which the E4 wristband can validly derive PR and PRV by comparing its UST interval recordings of mean PR and PRV metrics to that of a gold standard 5min interval recording. To achieve these two aims, participants' IBIs were simultaneously recorded with an E4 wristband and an ECG device during a 5min seated-rest condition.

#### 4.1 Validity of E4 wristband-based PRV

With regards to the E4 wristbands' PRV metrics validity as approximations of HRV metrics, our results are largely consistent with previous studies. Similar to Menghini et al. (2019), we found that, in a seated condition, the E4 wristband's estimation of mean PR, RMSSD, SDNN, and LF were comparable to their same mean HR and HRV estimation with the gold standard ECG device (see also van Lier et al., 2020). However, we found that HF was invalidly estimated by the E4 wristband. This finding is inconsistent with some results (e.g., Menghini et al., 2020; Schuurmans et al., 2020) but consistent with others (e.g., Ollander et al., 2016).

In line with Menghini et al. (2019), we found that the E4 wristband tended to overestimate RMSSD and LF relative to their estimation with the ECG device. This is valuable information for researchers as it allows researchers to potentially apply a calibration index (i.e., subtracting or adding a value) to the E4 wristband's estimation of PRV so that it approaches the true HRV value more precisely. However, the nature of this calibration index depends on whether the bias is proportional; as such, assessing proportional bias is needed<sup>2</sup>. Furthermore, it might be an interesting avenue for further inquiry to assess whether the bias (i.e., systematic error) of the E4 wristband is due to noise or some important underlying physiological parameter inherent to BVP. Analyzing the pulse transit time and pulse wave velocity may offer an opportunity to assess how, for instance, vasodilation/constriction, arterial stiffness, and arterial compliance (i.e., physiological factors influencing blood flow and pressure) are related to this systematic error (e.g., Mejia-Mejia et al., 2021; Mol et al., 2020).

#### 4.2 Validity of E4 wristband-based PRV obtained with UST intervals

The current study is, to our knowledge, the first to assess the fine-grained character of the time scales at which the E4 wristband can validly derive PR and PRV metrics. Our results corroborate many other studies (e.g., Baek et al., 2015; Munoz et al., 2015; Shaffer et al., 2016) in showing that mean PR/HR can be validly assessed at a large range of UST intervals (e.g., 10s). This indicates that the E4 wristband is, in that regard, similar to other recording devices. However, for the time-domain PRV metrics RMSSD and SDNN, the PRV recordings with the three shortest 10s UST intervals were unstable, whereas taking the average of them or using longer UST intervals (i.e., 30s and 120s) significantly increased stability. This result is not in line with Munoz et al. (2015), who demonstrated that RMSSD could be validly derived with 10s UST intervals (see also Nussinovitch et al., 2011). However, their participants were in a supine-rest condition as opposed to the seated-rest condition of the current study, perhaps making their recording more stable due to fewer motion artifacts. Other studies using a seated-rest condition also observed diminished stability of, for instance, RMSSD obtained with UST recording intervals of 10s (e.g., Baek et al., 2015; Salahuddin et al., 2007; Shaffer et al., 2016). Lu et al. (2008) compared PRV to HRV during supine and seated rest conditions and showed a diminishing accuracy during seated as compared to supine rest. The instability in the data might increase at such UST intervals of 10s because fluctuations in the PPI time series might occur precisely in the 10s recording, distorting the estimation of the time-domain PRV metrics RMSSD and SDNN. In longer recording intervals (e.g., 30s, 120s, 5min) or when averaging over multiple UST intervals, these irregular fluctuations in the PPI time series more likely are leveled out. Thus, averaging over multiple short intervals or increasing the UST recording length to 30s and 120s improved the

---

<sup>2</sup> For example, the E4 wristband's tendency to overestimate SDNN was negatively associated with the size of the ECG-based SDNN, implying that the calibration index shifts depending on ECG-based SDNN. Therefore, one would need the coefficients of the proportional bias model (i.e., calibration index =  $\beta_0 + \beta_1 \cdot \text{ECG-based value}$ ) to calculate the required calibration index.

accuracy of estimating RMSSD and SDNN with the E4 wristband. Especially for SDNN, this is not that surprising as it is a measure of the variability of the IBIs. This variability is expected to increase with the recording length (i.e., the more data, the more variability; Task Force, 1996). Our RMSSD and SDNN observations mostly mimic the results of other studies (e.g., Munoz et al., 2015; Shaffer et al., 2016). Although for the estimation with the UST interval of 30s, some found them to deviate too much from their 5min estimation (e.g., Shaffer et al., 2016), whereas others did find that these were reliable proxies (e.g., Baek et al., 2015; Munoz et al., 2015; Salahuddin et al., 2007). In general, for those UST intervals deemed valid in estimating mean PR and PRV metrics, as specified above, there was no tendency to over/underestimate mean PR or the PRV metrics relative to their same estimation with a 5min interval.

The HF obtained with the UST interval of 120s agreed insufficiently with its value acquired with 5min interval recording. This result is consistent with the finding of Shaffer et al. (2016), who found that HF could only be validly estimated with a UST interval of 180s. Even though others did find HF to be validly estimated with UST intervals of 20s and 40s (Baek et al., 2015; Salahuddin et al., 2007), we consider this implausible as a proper estimation of power in the HF band requires at least 70s (Task Force, 1996). This 70s minimum only ensures sufficient time to decently perform a power calculation in the HF band. This does not automatically imply that this is sufficiently long to estimate HF validly. As the current study examined BVP-based estimation of HF, which is less stable than ECG-based HF, even doubling this bare minimum 70s proved not to be sufficient. So based on our results, we recommend that longer BVP recordings are required for HF.

It is noteworthy in our second UST analysis, where we compared the estimation of the UST intervals of PR and PRV to the 5min ECG-based HR and HRV (see supplementary material D), that for RMSSD only the estimation with the UST interval of 30s was found valid. Therefore, the BVP-based estimation of RMSSD with UST intervals seems a less stable surrogate for ECG-based RMSSD with a 5min interval. On closer inspection of the Bland-Altman plots, we observed that, for the average of the three UST intervals of 10s and the UST interval of 120s, a large part (i.e., 92%) of these UST intervals' RMSSD estimation was acceptable, and only three observations mainly caused this invalidity with an E4 wristband-based inaccuracy of  $\pm 20$ ms. The discrepancy between the two UST analyses might be related to the BVP-related systematic error. As all measurements are BVP-based in the first analysis, this systematic error is accounted for. However, this is not the case in this second analysis. Here, systematic and random errors are combined, which might lead to differing results. The findings of the second UST analysis illustrate that BVP-based PRV mostly approximates ECG-based HRV. However, there are still marked differences between the two, likely due to different physiological constraints, as BVP-based PRV is bio-mechanical whereas ECG-based HRV is bio-electrical (see Yuda et al., 2020, for argumentation).

### 4.3 Statistical procedures

One source that leads to inconsistencies in research results and difficult comparability between studies is the statistical procedure employed to assess agreement between the proxy and the gold standard. Although commendable attempts have been made to establish valid statistical protocols for the assessment of agreement (e.g., Bland & Altman, 1999; Menghini et al., 2019; Menghini et al., 2021; Pecchia et al., 2018; van Lier et al., 2020), they still rely on some procedures that have been criticized such as correlational analyses and Cohen's *d* tests based on paired samples *t*-tests (e.g., Pecchia et al., 2018). For example, correlational analyses highlight the strength of an association but no consistent up/downward shift of the values of one of the variables relative to the other variable. We followed the proposed protocols to illustrate consistency in agreement patterns across several statistical procedures. We argue the Bland-Altman analysis should be the main source of information to identify agreement.

However, not all studies rely on Bland-Altman analysis (e.g., Ollander et al., 2016) or do so without using *a priori* LOAs and/or indicating how Bland-Altman assumptions were dealt with (e.g., Schuurmans et al., 2020; Kiran Kumar et al., 2021). As these assumptions strongly influence the calculation and representation of the 95% LOA, it is difficult to interpret the results of these studies unambiguously.

Moreover, the *a priori* LOA calculation is not consistent across the literature, and its proposed cut-off differs and is relatively arbitrary (i.e., 110% or 150%). Menghini et al. (2019) take  $\pm 50\%$  of the mean of the gold standard, whereas van Lier et al. (2020) take 10% of the range of biologically plausible HRV values, which can vary between age cohorts (e.g., 20-29 years or 50-60 years; Umetani et al., 1998). As highlighted in the methods section, we argue, similarly to Giavarina (2015), that this choice should depend on the goal for which the E4 wristband or UST interval estimations of PRV are used. For instance, for medical diagnostic purposes, the accuracy of PRV estimation as an approximation of HRV should be high (e.g., a priori LOA of 110%), whereas for lab-based research, it can be less strict (e.g., a priori LOA of 150%).

#### 4.4 Recommendations

E4 wristband validation studies examining PRV as approximations of HRV metrics and time scale at which PRV can validly be derived are limited or non-existing and use divergent procedures. Therefore, more procedurally consistent E4 wristband validation studies are needed.

Our study concludes that using the E4 wristband as a research-grade device to track PR/HR and estimate PRV/HRV in seated conditions in a lab-based context seems valuable under certain conditions. The E4 wristband provided comparable measures to the gold standard ECG device concerning mean HR, RMSSD, SDNN, and LF. This observation indicates that, in a lab setting where participants are seated and there is limited movement of the hand wearing the wristband, the E4 wristband can be a valid substitute for an ECG device with a 5min recording length. However, several observations need to be made.

First, the BVP signal was clearly less stable than the ECG signal (see supplementary material C), which led to a substantial exclusion of participants. Hence, the potential advantage of the E4 wristband to facilitate meeting sample size requirements more easily than with an ECG device might be canceled out by the sample size compensation needed to accommodate the instability of the BVP signal. However, as the E4 wristband is considerably cheaper, one can test multiple participants simultaneously. So even though one might need to exclude 32% of the data, it still may be more labor efficient to work with the E4 wristband than an ECG device. On the other hand, ethically, we must question whether designs, where 32% of the participants will have to be removed, are justified, especially in demanding task situations. We did not explicitly tell participants to keep the hand wearing the E4 wristband completely at rest during the baseline recording. It was only indicated to remain calm and control their emotions without excess movements (e.g., arm stretching). This perhaps caused an unstable BVP signal in some cases. On the other hand, as participants are likely to have made some minimal movement during the 5min rest, it shows that the E4 wristband's PR and PRV estimation as approximations of HR and HRV estimation was relatively valid under such minimal-movement conditions. We also asked participants to put on the E4 wristband themselves. Although we visually checked if it was put on correctly, it may have been that this was, on some occasions, not attached firmly enough, causing an unstable BVP signal. To enhance the quality of the signal, one could (1) ask participants to keep the hand wearing the E4 wristband completely at rest and (2) let the experimenter attach the E4 wristband watch to ensure that it is attached correctly and firmly. Furthermore, increasing the recording length might also remedy this. Longer BVP recording intervals lead to more stable BVP signal intervals that can be used to validly estimate the mean HR and HRV metrics (see, for example,

So et al., 2021). Lastly, it is important not to rely solely on BVP preprocessing algorithms (e.g., Kubios, Tarvainen et al., 2020) to correct artifacts and noise segments in the BVP signal. During preprocessing, we noticed that these algorithms sometimes missed uninterpretable noise segments. Therefore, we recommend a visual check of the BVP signal for its morphology and stability to ensure proper cleaning/preprocessing of the BVP signal. Moreover, it is advisable to include a rest period between trials of at least 2min (the current study used a 5min acclimatization period) to enhance the stationarity of the BVP data and, thus, its quality.

Second, we assessed the validity of the E4 wristband only in a condition involving minimal movement. But even then, we could not completely avoid the instability of the BVP signal. This observation possibly constrains the possibilities of using the E4 wristband in real-life situations. If movement is involved (e.g., walking or running), the E4 wristband is probably less likely to provide valid estimates. Still, other real-life situations might fulfill the E4 wristband's preconditions, such as studying HRV in classrooms or workplaces, at the bedside of patients, or at the home of older adults, where participants are seated or lying down and moving minimally. Furthermore, we only tested participants in a seated resting state. Thus, based on our results, we cannot assert how valid the E4 wristband is when, for instance, performing a computer task. However, future research-based studies can now implement the E4 wristband to assess its validity, for example, during task completion, under stress or when facing cognitive load. The recommendations we provide here can help to promote a valid E4 wristband data acquisition in a lab-based context. We also note that the current validation procedure was part of a broader creative problem solving experiment, of which the experimental setup might have had an undue influence on the validation procedure. However, Shcheslavskaya et al. (2010) have shown that the heart's activity returns to baseline within 6min after a cognitive challenge (see also Panaite et al., 2015). Although participants had to read instructions and perform six practice trials before the ECG and BVP recording, there was a 5min acclimatization period before the actual recording started, which should have been sufficient to eliminate any cardiovascular reactivity caused by the instructions and/or practice trials.

Third, our sample consisted almost exclusively of female undergraduates. Although this observation limits the scope in which our results can be interpreted, it is noteworthy that our results were similar to some other studies with a more biological sex-balanced sample (e.g., Menghini et al., 2019; van Lier et al., 2020) and having a broader age range (e.g., Menghini et al., 2019). Future studies should also consider recording additional demographic information, as it has been argued that, for instance, skin tone could impact PPG-based PRV recording (Fallow et al., 2013), although this has not been consistently reported (e.g., Bent et al., 2020).

Fourth, concerning the UST intervals of the E4 wristband, our results showed that the E4 wristband's time scales to derive PRV were not extremely fine-grained. Deriving PRV metrics based on 10s UST intervals did not produce valid estimates. However, the E4 wristband did prove its usefulness for lab experiments across all UST intervals for mean PR and when averaging over multiple short-duration intervals (e.g., 10s) or when using a recording length of at least 30s for RMSSD and SDNN. We must note that, solely for RMSSD, the UST intervals' comparison with 5min ECG only revealed a valid RMSSD with the UST interval of 30s. Although most RMSSD observations with the UST intervals were valid, a few deviated strongly from the 5min ECG-based RMSSD. Therefore, caution is warranted when using E4 wristband-based RMSSD with UST intervals as a proxy of ECG-based RMSSD. Assessing HF presumably would require a longer recording length than 120s (see, for example, Shaffer et al., 2016). These observations seem to place certain time constraints on using the E4 wristband in a lab context. For instance, studies aiming to assess PRV on a trial-by-trial basis with short trial durations (e.g., Shen et al., 2017) might not produce valid PRV estimates with the E4 wristband. However, the E4 wristband does lend itself to recording time-domain PRV metrics, for example, on a trial-by-trial basis or across a block of trials if the trial or block lasts at least 30s or if an



average can be taken across multiple shorter trials. For the frequency domain PRV metrics, a trial-by-trial analysis seems difficult as HF was found to be invalidly derived with the E4 wristband with a UST interval of 120s.

To summarize, the E4 wristband is a piece of promising research equipment in seated lab conditions. Our results showed that mean HR, RMSSD, SDNN, and LF were validly estimated by the E4 wristband. Contrary to some studies, we could not corroborate the valid estimation of HF. Furthermore, we showed that the E4 wristband's mean PR was a valid proxy of the 5min gold standard recording across all UST intervals. RMSSD and SDNN could be validly estimated with the E4 wristband with an average over multiple UST intervals of 10s, a UST interval of 30s or 120s, but not with 10s UST intervals. However, for RMSSD, when compared to a 5min ECG recording, only the UST interval of 30s remained valid. For HF, longer recording times seem to be required. Based on these results, we have formulated several recommendations for the use of the E4 wristband in laboratory research contexts.

### **FUNDING**

This work was supported by the "Fonds de la Recherche Scientifique" [grant number 34736358, 2019] and the Research Foundation Flanders [grant number G096919N].

### **ACKNOWLEDGMENTS**

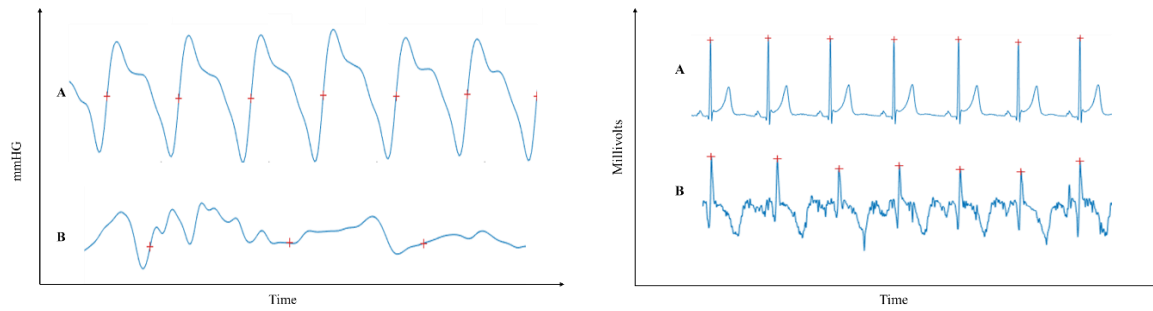
We thank Febe Demeyer, Christo Bratanov, Yujing Liang, and Amar Music for their assistance with the data collection and their critical thought and the Fonds de la Recherche Scientifique and the Research Foundation Flanders (FWO) for providing the opportunity to conduct this research under a research fellow grant.

## Appendix A. Supplementary material

### A. Normal and abnormal examples of BVP and ECG signals

**Figure 1**

*Normal and Abnormal examples of BVP and ECG signals*



*Note.* The left figure represents a normal (A) and abnormal (B) BVP signal morphology; the right figure represents a normal (A) and abnormal (B) ECG signal morphology; the Y-axis on the left figure represents the millimeter of mercury (mmHg), a unit of pressure; the Y-axis on the right figure represents millivolt.

## B. Data analysis with outliers included

### 1. Validity of E4 wristband-based PRV

#### 1.1 Step 1, PRV metric selection

Table 1 depicts the correlation coefficients and the descriptive statistics of the mean PR and HR, the time (RMSSD and SDNN), and frequency (LF and HF) domain PRV and HRV metrics. Correlations between both devices surpassing the cut-off of  $r = .70$  were found for HR, RMSSD, SDNN, LF, and HF. Therefore, the mean PR and all PRV metrics were included in step 2 of the three-step hierarchical procedure.

**Table 1.**

*Correlation Coefficients and the Descriptive Statistics of PR/HR and PRV/HRV Metrics obtained with E4 and ECG*

Metrics	$r$ [95% CI]	$M(SD)$ E4	$M(SD)$ ECG
PR/HR	.9989[.9983, .9996]	81.75(7.41)	81.90(7.34)
RMSSD	.95[.92, .99]	41.40(14.33)	35.52(15.00)
SDNN	.987[.978, .998]	46.57(15.91)	45.03(16.52)
LF	.95[.90, .97]	1126.52(797.46)	1159.06(814.12)
HF	.98[.95, 1.00]	886.24(768.88)	742.43(654.03)

*Note.* PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); LF, low frequency (ms<sup>2</sup>); HF, high frequency (ms<sup>2</sup>);  $r$ , Pearson correlation coefficient; 95% CI, 95% confidence intervals.

#### 1.2 Step 2, PRV metric validity

Table 2 presents the results, and Figure 1 presents the Bland-Altman plots (see supplementary material E for the assumptions handling). Concerning mean PR/HR, there was a slight tendency of the E4 wristband to underestimate it relative to its estimation by the ECG device. Namely, the upper bound of the bias's 95% CI (i.e., -0.006) was just below zero. However, the 95% LOA was within the *a priori* 150% LOA, with 100% of the observed differences lying within the *a priori* 150% LOA bounds. These findings illustrated that the deviations between the mean PR obtained with the E4 wristband and the mean HR acquired with the ECG device were within the maximum acceptable deviation limits.

For RMSSD and SDNN, the 95% CI of its bias was entirely above zero, indicating that the E4 wristband tended to overestimate RMSSD and SDNN relative to the same metric acquired with the ECG device. For RMSSD and SDNN, the 95% LOA was within the *a priori* 150% LOA, with 100% and 100% of the observed differences lying within the *a priori* 150% LOA bounds, respectively. Therefore, the E4 wristband's estimation of RMSSD and SDNN may be considered sufficiently in agreement with the same metrics estimated by the ECG device.

For LF and HF, proportional bias was detected (i.e., a positive association between the observed differences and the ECG device estimation of LF and HF; see Figure 1). LF was overestimated by the E4 wristband for observations of ECG-based LF higher than 436ms<sup>2</sup>, whereas the bias was non-significant for lower ECG-based LF values. And, HF was overestimated by the E4 wristband for observations of ECG-based HF between 47ms<sup>2</sup> and 1719ms<sup>2</sup>, whereas the bias was non-significant for ECG-based HF values falling outside that range. For LF, the 95% LOA were within the *a priori* 150%

LOA for 85% of the ECG-based LF observations, with acceptable 95% LOA for observations of ECG-based LF lower than 2168ms<sup>2</sup>. Here, 95% of the observed differences were lying within the *a priori* 150% LOA bounds. As such, the E4 wristband estimation of LF was in sufficient agreement with its estimation by the ECG device. However, for HF, the 95% LOA were within the *a priori* 150% LOA for 46% of the ECG-based HF observations, with acceptable 95% LOA for observations of ECG-based HF lower than 521ms<sup>2</sup>. Here, only 90% of the observed differences were within the *a priori* 150% LOA. This finding illustrated that the E4 wristband's estimation of HF deviated too much from its gold standard estimation by the ECG device.

Consequently, the mean PR, RMSSD, SDNN, and LF were retained for the last step of the three-step hierarchical procedure.

**Table 2.**

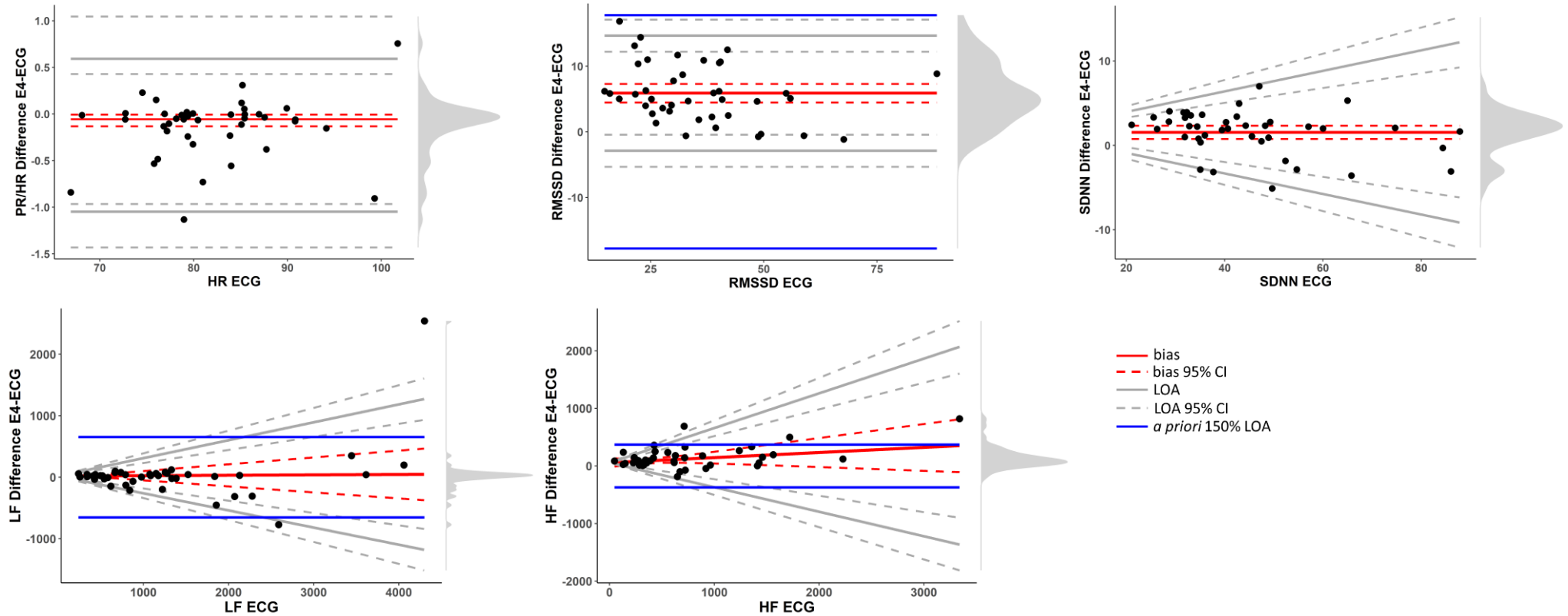
*Bland-Altman Analysis: Bias, 95% LOA, and A Priori 150% LOA*

Metrics	Bias [95% CI]	Lower LOA [95% CI]	Upper LOA [95% CI]	<i>A priori</i> 150% LOA
PR/HR*	-0.06 [-0.13, -0.006]	-1.05 [-1.43, -0.97]	0.59 [0.43, 1.05]	-41.00, 41.00
RMSSD*	5.86 [4.74, 7.30]	-2.88 [-5.32, -0.43]	14.65 [12.21, 17.10]	-17.76, 17.76
SDNN*	1.53 [0.76, 2.33]	Bias - 0.12*GS [0.09, 0.16]	Bias + 0.12*GS [0.09, 0.16]	-22.52, 22.52
LF*	-18.36 + 0.006*GS [-50.91, 22.53], [-0.01, 0.06]	Bias - 0.29*GS [0.21, 0.36]	Bias + 0.29*GS [0.21, 0.36]	-653.21, 653.21
HF	55.92 + 0.09*GS [-23.99, 82.11], [0.01, 0.24]	Bias - 0.51*GS [0.38, 0.65]	Bias + 0.51*GS [0.38, 0.65]	<b>-371.21, 371.21</b>

*Note.* PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); LF, low frequency (ms<sup>2</sup>), HF, high frequency (ms<sup>2</sup>); *PR/HR distribution characteristics*, homoscedastic and non-normal which a log transform could not alleviate; *RMSSD distribution characteristics*, homoscedastic and normal; *SDNN distribution characteristics*, heteroscedastic and non-normal which a log transform could alleviate; *LF distribution characteristics*, heteroscedastic and non-normal which log transformation could not alleviate; *HF distribution characteristics*, heteroscedastic and non-normal which a log transform could not alleviate (see supplementary material E for the assumption handling); Median; 2.5 percentile; 97.5 percentile;  $\beta_0$ , the intercept;  $\beta_1$ , the slope coefficient; antilog, the antilog slope value; GS, the gold standard value; LOA, limits of agreement; 95% CI, 95% confidence intervals; *a priori* 150% LOA, *a priori* 150% limits of agreement; bold typeface indicates when the 95% LOA is outside the *a priori* 150% LOA; \* indicates the proxies that were retained for the following step in the three-step hierarchical procedure.

Figure 1

Bland-Altman Plots for Mean PR/HR and PRV/HRV Metrics



Note. PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); LF, low frequency ( $\text{ms}^2$ ), HF, high frequency ( $\text{ms}^2$ ); density distribution of the differences at the right side of each plot; 95% CI, 95% confidence interval; LOA, limits of agreement; *a priori* LOA of 150% (only displayed if close to 95% LOA).

### *1.3 Step 3, the magnitude of difference*

The observed effects sizes were negligible for mean PR/HR (Cliff's  $d = -.02$ , 95% CI [-.27, .23]). For the time-domain metrics, the effect sizes were small for RMSSD (Cliff's  $d = .27$ , 95% CI [.02, .49]) and negligible for SDNN (Cliff's  $d = .09$ , 95% CI [-.16, .34]). For the frequency-domain metric LF, the effect size was negligible (Cliff's  $d = -.005$ , 95% CI [-.25, .24]). These findings show that the E4 wristband provided comparable estimates of mean PR, RMSSD, SDNN, and LF to their same mean HR and HRV estimation with the ECG device.

## **2 Validity of E4 wristband-based PRV obtained with UST intervals**

### *2.1 E4 wristband UST interval PRV versus 5min E4 wristband-based PRV*

**2.1.1 Step 1, PRV metric selection.** The correlation coefficients and descriptive statistics of the mean PR and the PRV metrics of the UST intervals are presented in Table 3. Correlations between UST intervals and 5min intervals surpassing the cut-off of  $r = .70$  were found for all UST intervals estimating PR, RMSSD, SDNN, and HF, except for one of the three UST intervals of 10s for RMSSD and SDNN. Hence, all UST intervals estimating PR and PRV surpassing the cut-off were included in the second step of the three-step hierarchical procedure.

**Table 3.**

*Correlation Coefficients and Descriptive Statistics of the PR and PRV Metrics for the UST and 5min Recordings Obtained with the E4 Wristband*

Metrics	<i>r</i> (95% CI)	<i>M</i> ( <i>SD</i> )
<i>PR</i>		
5min		81.75(7.41)
10s T1	.91[.86, .98]	81.96(8.07)
10s T2	.89[.82, .98]	81.71(8.00)
10s T3	.94[.89, .99]	80.94(8.76)
Average of 10s	.96[.94, .99]	81.53(7.86)
30s	.91[.84, .99]	82.01(7.60)
120s	.98[.96, 1.00]	81.68(7.07)
<i>RMSSD</i>		
5min		41.41(14.33)
10s T1	.58[.35, .80]	42.12(19.20)
10s T2	.82[.68, 1.00]	39.51(21.03)
10s T3	.79[.64, .97]	38.26(17.35)
Average of 10s	.90[.77, 1.00]	39.96(15.69)
30s	.89[.79, 1.00]	40.22(16.76)
120s	.96[.93, 1.00]	42.07(15.54)
<i>SDNN</i>		
5min		46.56(15.92)
10s T1	.58[.34, .84]	51.66(23.91)
10s T2	.72[.55, .89]	41.55(22.48)
10s T3	.79[.66, .94]	41.49(18.66)
Average of 10s	.86[.74, 1.00]	44.90(17.29)
30s	.86[.76, .98]	46.13(20.93)
120s	.95[.91, 1.00]	47.62(16.80)
<i>HF</i>		
5min		886.25(768.88)
120s	.92[.85, 1.00]	857.56(778.45)

*Note.* PR, mean pulse rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); T, time interval; *r*, Pearson correlation coefficient; 95% CI = 95% confidence interval.

**2.1.2 Step 2, PRV metric validity.** The results of the Bland-Altman analysis are depicted in Table 4, and exemplary Bland-Altman plots are shown in Figure 2 (see supplementary material E for the assumptions handling). The bias included zero for all UST intervals estimating mean PR in its 95% CI, except for the UST interval of 120s. Here, proportional bias was detected (i.e., a negative association between the observed differences and the 5min interval estimation of mean PR; see Figure 2). Specifically, PR was slightly underestimated by the UST interval of 120s for observations of 5min-based PR higher than 99 bpm, whereas the bias was non-significant for lower 5min-based PR values. For all other UST intervals, no tendency to over/underestimate mean PR was found relative to its same estimation with the 5min interval. For all UST intervals' estimation of mean PR, the 95% LOA was

inside the *a priori* 150% LOA, thereby showing that the UST intervals provided comparable mean PR values to the same value acquired with the 5min interval. In all cases, 100% of the data lay within the *a priori* 150% LOA.

For all the UST intervals estimating RMSSD except for one of the UST intervals of 10s (i.e., T3), the bias included zero in its 95% CI. Therefore, for all these UST intervals, no tendency to over/underestimate RMSSD relative to the same estimation with the 5min interval was found. For the UST interval of 10s (i.e., T3), the 95% CI of the bias was entirely below zero, indicating a tendency of the UST interval of 10s to underestimate RMSSD relative to its same estimation with 5min interval. For the UST interval of 10s (i.e., T3), the 95% LOA was within the *a priori* 150% LOA for 34% of the 5min-based RMSSD observations, with acceptable 95% LOA for observations of 5min-based RMSSD lower than 34ms (see Figure 2). Here, 95% of the observed differences were lying within the *a priori* 150% LOA bounds. For the UST interval of 30s, the 95% LOA was within the *a priori* 150% LOA for 88% of the 5min-based RMSSD observations, with acceptable 95% LOA for observations of 5min-based RMSSD lower than 56ms (see Figure 2). Here, 98% of the observed differences were lying within the *a priori* 150% LOA bounds. For the UST intervals of 120s and the average of the three UST intervals of 10s, the 95% LOA were inside the *a priori* 150%, with 98% and 100% of the observed differences inside the *a priori* 150% LOA boundaries, respectively. Therefore, the RMSSD estimated with the UST interval of 10s (i.e., T3), the average of the three UST intervals of 10s, the UST intervals of 30s and 120s were in sufficient agreement with its estimation by the 5min interval. However, for the UST intervals of 10s (i.e., T2) the 95% LOA was within the *a priori* 150% LOA for 17% of the 5min-based RMSSD observations, with acceptable 95% LOA for observations of 5min-based RMSSD lower than 29ms. Here, only 88% of the observed differences were lying within the *a priori* 150% LOA bounds. Therefore, the estimation of RMSSD with this UST interval of 10s (i.e., T2) deviated too much from the same value acquired with the 5min interval, indicating insufficient agreement.

For all UST intervals estimating SDNN, the bias included zero in its 95% CI, except for the two UST intervals of 10s. Here, the upper bound of the 95% CI was below zero, indicating a tendency of these UST interval to underestimate SDNN relative to its estimation with the 5min interval. For all other UST intervals estimating SDNN, no tendency to over/underestimate it was found relative to its same estimation with the 5min interval. The 95% LOA was within the *a priori* 150% LOA for the average of the three UST intervals of 10s and the UST interval of 120s. For these UST intervals, 98% and 100% of the observed differences were inside the *a priori* 150% LOA boundaries, respectively. For the UST interval of 10s (i.e., T3) estimating SDNN, the lower bound of the 95% LOA was outside the *a priori* 150% LOA, with 95% of the observed differences within the *a priori* 150% LOA boundaries. For the UST interval of 30s estimating SDNN, the 95% LOA was inside the *a priori* 150% LOA for the 5min-based SDNN observations below 49ms, including 66% of the SDNN observations (see Figure 2). 98% of the observed differences were inside the *a priori* 150% LOA boundaries. This result showed that the SDNN estimated by the UST interval of 10s (i.e., T3), the average of the three UST intervals of 10s, the UST interval of 30s and 120s were in sufficient agreement with its same estimation by the 5min interval. For the UST intervals of 10s (i.e., T2) estimating SDNN, the lower bound of the 95% LOA was outside *a priori* 150% LOA, with 88% of the observed differences inside the *a priori* 150% LOA. As such, the SDNN estimated with the UST interval of 10s (i.e., T2) deviated too much from its estimation with the 5min interval, thereby showing insufficient agreement.

For the UST interval of 120s estimating HF, the upper bound of the bias's 95% CI was below zero, indicating a tendency of this UST interval to underestimate HF relative to its same estimation with the 5min interval (see Figure 3). Furthermore, the 95% LOA were within the *a priori* 150% LOA for 5% of the 5min-based HF observations, with acceptable 95% LOA for observations of 5min-based HF lower than 178ms<sup>2</sup>. Here, only 85% of the observed differences were within the *a priori* 150% LOA



boundaries. Therefore, the HF estimated with the UST interval of 120s deviated too much from its estimation with the 5min interval, indicating insufficient agreement.

Consequently, all UST intervals estimating mean PR, one of the three UST intervals of 10s (i.e., T3), the average of the three UST intervals of 10s, the UST interval of 30s and 120s estimating RMSSD and SDNN, were further analyzed in the final step of the three-step hierarchical procedure.

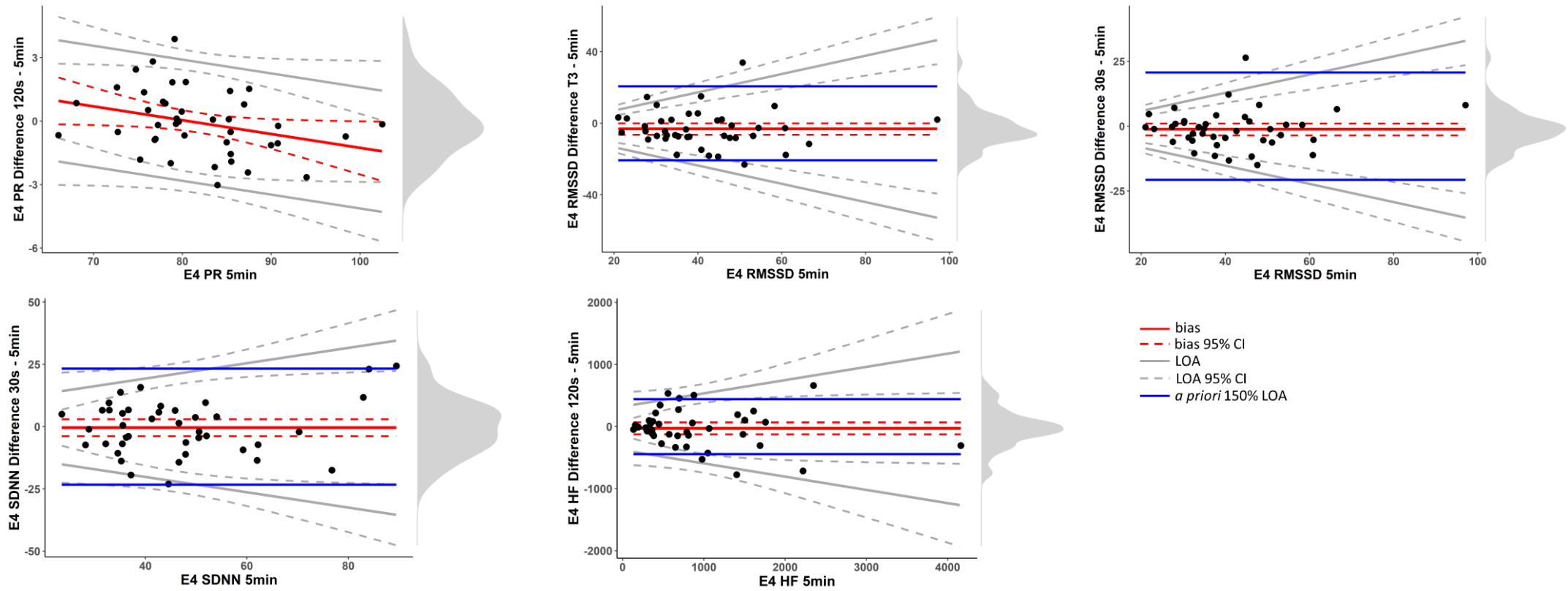
**Table 4.***UST Recording's Bland-Altman Analysis: Bias and 95% LOA*

Metrics	Bias [95% CI]	Lower LOA [95% CI]	Upper LOA [95% CI]	<i>A priori</i> 150% LOA
<i>PR</i>				
10s T1*	0.20 [-0.84, 1.24]	-6.42 [-8.04, -4.45]	6.65 [4.85, 8.45]	-40.88, 40.88
10s T2*	-0.04 [-1.18, 1.10]	-7.12 [-9.09, -5.14]	7.04 [5.06, 9.01]	-40.88, 40.88
10s T3*	-0.82 [-1.82, 0.18]	-7.04 [-8.77, -5.30]	5.40 [3.66, 7.13]	-40.88, 40.88
Average of 10s*	-0.22 [-0.89, 0.45]	-4.37 [-5.53, -3.21]	3.93 [2.77, 5.09]	-40.88, 40.88
30s*	0.26 [-0.75, 1.27]	-6.00 [-7.75, -4.26]	6.52 [4.77, 8.26]	-40.88, 40.88
120s*	<b>5.27</b> - 0.07*GS [0.03, 10.52], [-0.13, 0.00]	Bias - 1.96*1.46 [1.21, 1.79]	Bias + 1.96*1.46 [1.21, 1.79]	-40.88, 40.88
<i>RMSSD</i>				
10s T2	-1.90 [-5.78, 1.98]	Bias - 2.46*(-2.49 + 0.28*GS) [-9.16, 4.17], [0.12, 0.43]	Bias + 2.46*(-2.49 + 0.28*GS) [-9.16, 4.17], [0.12, 0.43]	<b>-20.70, 20.70</b>
10s T3*	-3.14 [-6.44, -0.03]	Bias - 0.51*GS [0.37, 0.65]	Bias + 0.51*GS [0.37, 0.65]	<b>-20.70, 20.70</b>
Average of 10s*	<b>-1.17</b> [-4.64, 1.58]	<b>-20.55</b> [-33.41, -15.51]	<b>14.75</b> [11.05, 23.78]	-20.70, 20.70
30s*	-1.18 [-3.52, 1.06]	Bias - 0.35*GS [0.25, 0.45]	Bias + 0.35*GS [0.25, 0.45]	-20.70, 20.70
120s*	0.67 [-0.70, 2.04]	Bias - 2.46*(0.43 + 0.07*GS) [-1.96, 2.81], [0.02, 0.13]	Bias + 2.46*(0.43 + 0.07*GS) [-1.96, 2.81], [0.02, 0.13]	-20.70, 20.70
<i>SDNN</i>				
10s T2	-5.01 [-9.95, -0.07]	Bias - 2.46*(2.57 + 0.20*GS) [-7.02, 12.16], [0.003, 0.39]	Bias + 2.46*(2.57 + 0.20*GS) [-7.02, 12.16], [0.003, 0.39]	<b>-23.28, 23.28</b>
10s T3*	-5.07 [-8.73, -1.42]	-27.77 [-34.10, -21.44]	17.62 [11.29, 23.96]	<b>-23.28, 23.28</b>
Average of 10s*	-1.66 [-4.46, 1.13]	-19.00 [-23.83, -14.16]	15.67 [10.83, 20.50]	-23.28, 23.28
30s*	-0.43 [-3.85, 2.99]	Bias - 2.46*(3.03 + 0.13*GS) [-2.31, 8.37], [0.02, 0.23]	Bias + 2.46*(3.03 + 0.13*GS) [-2.31, 8.37], [0.02, 0.23]	-23.28, 23.28
120s*	1.05 [-0.63, 2.74]	-9.41 [-12.33, -6.49]	11.51 [8.59, 14.43]	-23.28, 23.28
<i>HF</i>				
120s	-28.69 [-125.00, -67.70]	Bias - 2.46*(143 + 0.09*GS) [48.50, 237.00], [0.006, 0.17]	Bias + 2.46*(143 + 0.09*GS) [48.50, 237.00], [0.006, 0.17]	<b>-443.12, 443.12</b>

Note. PR, mean pulse rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); T = time interval; *PR UST intervals distribution characteristics*, all were homoscedastic and normally distributed with only the UST interval of 120s showing proportional bias; *RMSSD UST distribution characteristics*, all were homoscedastic except for T2 and the UST interval of 120s. T3, the average of the three UST interval of 10s, and the UST interval of 30s displayed non-normality for which log transformation only alleviated non-normality for T3 and the UST interval of 30s; *SDNN UST intervals distribution characteristics*, all were homoscedastic and normally distributed except for T2 and the UST interval of 30s which were heteroscedastic; *HF UST interval distribution characteristics*, data was heteroscedastic and normally distributed (see supplementary material E for the assumption handling); Median; 2.5 percentile; 97.5 percentile;  $\beta_0$ , the intercept;  $\beta_1$ , the slope coefficient; Bias SD, the SD of the residuals of the proportional bias model; antilog, the antilog slope value; GS, the gold standard value; LOA = 95% limits of agreement; 95% CI = 95% confidence interval; *A priori* 150% LOA, *a priori* 150% limits of agreement; Bold typeface = the 95% LOA is outside the *a priori* 150% LOA; \* indicates the proxies that were retained for the following step in the three-step hierarchical procedure.

**Figure 2**

*Examples of Bland-Altman Plots Comparing Mean PR and PRV Metrics Obtained with UST Intervals by the E4 Wristband Versus Their 5min Recording with the E4 Wristband*



*Note.* PR, mean pulse rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); AVG10s, average of the three 10s intervals; T3, third UST interval of 10s; density distribution of the differences at right side plot; 95% LOA, 95% limit of agreement; 95% CI, 95% confidence intervals; the *a priori* LOA of 150% only displayed if close to 95% LOA.

**2.1.3 Step 3, magnitude of difference.** Regarding mean PR, the effect sizes were negligible for all the UST intervals: T1 (Cliff's  $d = -.02$ , 95% CI [-.27, .22]), T2 (Cliff's  $d = .03$ , 95% CI [-.22, .27]), T3 (Cliff's  $d = .11$ , 95% CI [-.14, .34]), the average of three 10s intervals (Cliff's  $d = .03$ , 95% CI [-.22, .27]), 30s interval (Cliff's  $d = -.02$ , 95% CI [-.26, .23]) and 120s interval (Cliff's  $d = -.0006$ , 95% CI [-.25, .25]). Concerning RMSSD, a small effect sizes was found for UST interval of 10s (Cliff's  $d = .19$ , 95% CI [-.06, .42]) and negligible effect sizes were found for the average of the three UST intervals of 10s (Cliff's  $d = .09$ , 95% CI [-.16, .33]), the UST interval of 30s (Cliff's  $d = .11$ , 95% CI [-.14, .35]), and the UST interval of 120s (Cliff's  $d = -.04$ , 95% CI [-.28, .21]). Lastly, for the UST intervals estimating SDNN a small effect size was found for the UST interval of 10s (Cliff's  $d = .25$ , 95% CI [-.006, .47]) and negligible effect sizes were found for the average of the three UST intervals of 10s (Cliff's  $d = .09$ , 95% CI [-.16, .33]), the UST interval of 30s (Cliff's  $d = .03$ , 95% CI [-.21, .26]), and the UST interval of 120s (Cliff's  $d = -.03$ , 95% CI [-.28, .22]). These results show that the UST intervals specified above provide estimates of mean PR, RMSSD and SDNN comparable to the same value estimated with a 5min interval.

## 2.2 E4 wristband UST interval PRV versus 5min ECG-based HRV

**2.2.1 Step 1, PRV metric selection.** For all UST intervals, strong correlations ( $r > .70$ ) were observed between their estimation of PR and PRV and the 5min ECG-based estimation of HR and HRV, except for one of the three UST intervals of 10s for RMSSD and SDNN.

**2.2.2 Step 2, PRV metric validity.** Considering only those UST intervals that survived step 1, all UST intervals' estimations of PR were in sufficient agreement with the 5min ECG-based estimation of HR. Concerning RMSSD, only its estimation with the UST interval of 30s was in sufficient agreement with the same estimation with the 5min ECG recording. For the two UST interval of 10s, the average of the three UST intervals of 10s and the UST interval of 120s, only 83%, 93%, 93% and 90% of the observed differences were within the *a priori* 150% LOA, respectively. Regarding SDNN, only its estimation with the average of the three UST intervals of 10s and the UST interval of 30s and 120s were in sufficient agreement with their same estimation with the 5min ECG recording. For the SDNN estimation with the two UST interval of 10s, only 83% and 93% of the observed differences were within the *a priori* 150% LOA, respectively. Concerning HF, its estimation with the UST interval of 120s agreed insufficiently with the same estimation with the 5min ECG recording.

**2.2.3 Step 3, magnitude of difference.** Considering only those UST intervals that survived step 2, negligible effect sizes were found between the UST intervals estimation of PR and SDNN and the 5min ECG-based estimation of HR and SDNN. A small effect size was found for the UST interval of 30s estimating RMSSD.

### C. The percentage of corrected IBI and the percentage of noise free ECG/BVP recording

**Table 1**

*Statistics Representing the Percentage of Corrected IBI and the Percentage of Noise Free ECG/BCP Recording Interval*

	BVP				ECG			
	%IBI corrected		%noise free		%IBI corrected		%noise free	
	<i>M(SD)</i>	range	<i>M(SD)</i>	range	<i>M(SD)</i>	range	<i>M(SD)</i>	range
5min	1.28(1.24)	0-4.37	95.96(4.41)	83.95-100	0.238(0.53)	0-2.77	99.92(0.32)	98.33-100
30s	0.58(1.40)	0-4.76	99.38(3.11)	80.00-100				
120s	1.13(1.38)	0-4.88	97.55(3.95)	83.33-100				

*Note.* BVP, blood-volume pulse; ECG, electrocardiogram; %IBI, percentage of corrected inter-beat-intervals of the noise free segment, %noise free; percentage of the recording interval that was noise free and used to calculate the HRV metrics; 5min, baseline recording interval; 30s, the E4 wristband ultra-short BVP recording interval of 30s; 120s, the E4 wristband ultra-short BVP recording interval of 30s. ECG was only recorded for a 5min baseline interval.

## **D. E4 wristband UST interval PRV versus 5min ECG-based HRV**

### **1. Step 1, PRV metric selection**

The correlation coefficients and descriptive statistics of the mean PR/HR and the PRV/HRV metrics of the UST intervals are presented in Table 1. For mean PR, all UST intervals illustrated strong correlations with the mean HR acquired with the ECG device's 5min interval. For the time-domain PRV metrics RMSSD and SDNN, strong correlations were found between their estimation with UST intervals and their estimation with ECG device's 5min interval, except for the three UST intervals of 10s for RMSSD and two of the three UST intervals of 10s for SDNN. Lastly, for the frequency-domain PRV metric HF, a strong correlation was observed between its estimation with the UST interval of 120s and its estimation with the ECG device's 5min interval. As the intervals mentioned above were the sole ones that surpassed the a priori defined cut-off of  $r = .70$ , they were included in the second step of the three-step hierarchical procedure.

**Table 1.**

*Correlation Coefficients and Descriptive Statistics of the PR and PRV Metrics for the UST Obtained with the E4 Wristband and the HR and HRV metrics acquired with the 5min ECG Recording*

metrics	<i>r</i> (95% CI)	<i>M</i> ( <i>SD</i> )
<i>PR/HR</i>		
5min		82.20(7.10)
10s T1	.90[.84, .98]	82.22(7.80)
10s T2	.89[.81, .98]	82.00(8.00)
10s T3	.93[.88, .98]	81.27(8.51)
Average of 10s	.96[.93, .99]	81.83(7.67)
30s	.90[.83, .99]	82.19(7.40)
120s	.98[.96, 1.00]	82.05(6.74)
<i>RMSSD</i>		
5min		33.34(11.43)
10s T1	.45[.14, .77]	39.64(15.55)
10s T2	.61[.42, .81]	36.78(14.95)
10s T3	.66[.48, .83]	36.27(14.42)
Average of 10s	.76[.57, .95]	37.56(11.26)
30s	.74[.60, .87]	37.72(12.23)
120s	.85[.77, .94]	39.80(10.94)
<i>SDNN</i>		
5min		42.92(13.91)
10s T1	.50[.23, .80]	49.24(20.85)
10s T2	.59[.36, .81]	38.84(18.22)
10s T3	.73[.52, .99]	39.93(17.59)
Average of 10s	.78[.59, 1.00]	42.67(14.47)
30s	.74[.57, .95]	42.83(15.20)
120s	.92[.88, .98]	45.59(14.44)
<i>HF</i>		
5min		637.72(450.89)
120s	.88[.79, .97]	725.48(511.31)

*Note.* PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); T, time interval; *r*, Pearson correlation coefficient; 95% CI = 95% confidence interval.

## 2. Step 2, PRV metric validity

The results of the Bland-Altman analysis are depicted in Table 2, and exemplary Bland-Altman plots are shown in Figure 1 (see supplementary material E for the assumptions handling). The bias included zero for all UST intervals estimating mean PR in its 95% CI, except for the for the UST interval of 120s. Here, proportional bias was detected (i.e., a negative association between the observed differences and the 5min interval estimation of mean PR; see Figure 1). Specifically, PR was slightly underestimated by the E4 wristband for observations of ECG-based HR higher than 91 bpm, whereas

the bias was non-significant for lower ECG-based HR values. For all other UST intervals, no tendency to over/underestimate mean PR was found relative to its same estimation with the 5min interval. For all UST intervals' estimation of mean PR, the 95% LOA was inside the *a priori* 150% LOA, thereby showing that the UST intervals provided comparable mean PR values to the same value acquired with the 5min interval. In all cases, 100% of the data lay within the *a priori* 150% LOA.

For the UST interval of 30s estimating RMSSD, the bias included zero in its 95% CI (see Figure 1). Therefore, there was no tendency to over/underestimate RMSSD relative to the same estimation with the ECG device's 5min interval for this UST interval. The upper bound of the 95% LOA was outside the *a priori* 150% LOA, with 95% of the observed differences inside the *a priori* 150% LOA. This finding illustrated that the RMSSD estimated with the UST intervals of 30s was borderline in sufficient agreement with its estimation by the ECG device's 5min interval. For the average of the three UST intervals of 10s and the UST interval of 120s, proportional bias was detected (see Figure 1). In both cases, a negative association was found between the observed differences and the ECG device's 5min interval estimation of RMSSD. Specifically, RMSSD was overestimated by the average of the three UST intervals of 10s and the UST interval of 120s for observations for ECG-based RMSSD lower than 44ms and 50ms, respectively. In contrast, the bias was nonsignificant for ECG-based RMSSD values higher than 44ms and 50ms. Furthermore, for the average of the three UST intervals of 10s and the UST interval of 120s, the 95% LOA were inside the *a priori* 150% LOA for 17% and 36% of the ECG-based RMSSD observations, respectively, with acceptable 95% LOA for observations of ECG-based RMSSD higher than 46ms and 39ms, respectively. In both UST interval cases, only 92% of the observed differences were inside the *a priori* 150% LOA. Therefore, the estimation of RMSSD with the average of the three UST intervals of 10s and the UST interval of 120s deviated too much from the same value acquired with the ECG device's 5min interval, indicating insufficient agreement.

For all UST intervals estimating SDNN, the bias included zero in its 95% CI, except for the UST intervals of 120s. Here, the lower bound of the 95% CI was above zero, indicating a tendency of this UST interval to overestimate SDNN relative to its estimation with the ECG device's 5min interval. For all other UST intervals estimating SDNN, no tendency to over/underestimate it was found relative to its same estimation with the ECG device's 5min interval. The 95% LOA was within the *a priori* 150% LOA for the UST interval of 30s, and 120s. For these UST intervals, 97% and 100% of the observed differences were inside the *a priori* 150% LOA boundaries, respectively. For the average of the three UST intervals of 10s, the lower bound of the 95% LOA was outside the *a priori* 150% LOA (see Figure 1). Nonetheless, 97% of the observed differences were still within the *a priori* 150% LOA boundaries. This result showed that the SDNN estimated by the average of the three UST intervals of 10s, the UST interval of 30s, and 120s were in sufficient agreement with its same estimation by the ECG device's 5min interval. For the UST interval of 10s, the 95% LOA was inside the *a priori* 150% LOA for 23% of the ECG-based SDNN observations, with acceptable 95% LOAs for observations of ECG-based SDNN lower than 32ms. Only 90% of the observed differences were within the *a priori* 150% LOA boundaries. As such, the SDNN estimated with the UST interval of 10s deviated too much from its estimation with the 5min interval, thereby showing insufficient agreement.

For the UST interval of 120s estimating HF, the lower bound of the bias's 95% CI was above zero, indicating a tendency of this UST interval to overestimate HF relative to its same estimation with the ECG device's 5min interval (see Figure 1). Furthermore, the 95% LOA was outside the *a priori* 150% LOA, with only 77% of the observed differences within the *a priori* 150% LOA boundaries. Therefore, the HF estimated with the UST interval of 120s deviated too much from its estimation with the 5min interval, indicating insufficient agreement.

Consequently, all UST intervals estimating mean PR, the UST interval of 30s estimating RMSSD, and the average of the three UST intervals of 10s, the UST interval of 30s and 120s estimating SDNN, were further analyzed in the final step of the three-step hierarchical procedure.



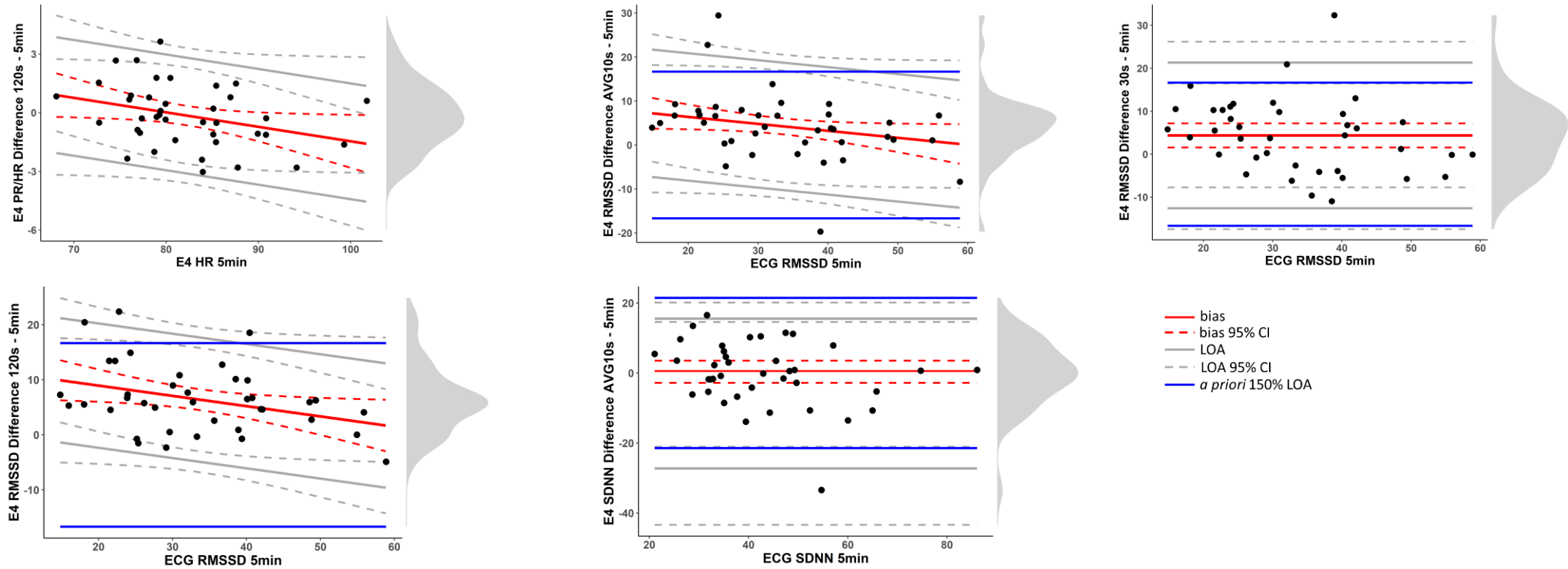
**Table 2.***Bland-Altman Analysis for UST recordings Versus ECG 5min recording: Bias and 95% LOA*

metric	Bias [95% CI]	Lower LOA [95% CI]	Upper LOA [95% CI]	<i>A priori</i> 150% LOA
<b>PR</b>				
10s T1*	0.03 [-1.06, 1.12]	-6.56 [-8.44, -4.67]	6.61 [4.73, 8.50]	-41.10, 41.10
10s T2*	-0.19 [-1.37, 0.98]	-7.31[-9.35, -5.27]	6.92 [4.88, 8.96]	-41.10, 41.10
10s T3*	-0.93 [-1.98, 0.13]	-7.31[-9.14, -5.48]	5.46 [3.63, 7.29]	-41.10, 41.10
Average of 10s*	-0.36 [-1.08, 0.35]	-4.67[-5.91, -3.44]	3.95 [2.71, 5.18]	-41.10, 41.10
30s*	-0.003 [-1.05, 1.04]	-6.31[-8.12, -4.50]	6.30 [4.50, 8.11]	-41.10, 41.10
120s*	5.94 - 0.07*GS [0.09, 11.79], [-0.14, 0.00]	Bias - 1.96*1.51 [1.26, 1.86]	Bias + 1.96*1.51 [1.26, 1.86]	-41.10, 41.10
<b>RMSSD</b>				
Average of 10s	9.56 - 0.15*GS [6.12, 12.61], [-0.29, -0.04]	Bias - 1.96*7.13 [5.04, 10.45]	Bias + 1.96*7.13 [5.04, 10.45]	-16.67, <b>16.67</b>
30s*	4.37 [1.57, 7.18]	-12.58 [-17.44, -7.73]	21.33 [16.47, 26.18]	-16.67, <b>16.67</b>
120s	12.70 - 0.19*GS [6.77, 18.59], [-0.35, -0.02]	Bias - 1.96*5.76 [4.65, 7.28]	Bias + 1.96*5.76 [4.65, 7.28]	-16.67, <b>16.67</b>
<b>SDNN</b>				
10s T3	-2.99 [-6.91, 0.93]	Bias - 2.46*(0.76 + 0.21*GS) [-5.90, 7.42], [0.06, 0.36]	Bias + 2.46*(0.76 + 0.21*GS) [-5.90, 7.42], [0.06, 0.36]	<b>-21.46, 21.46</b>
Average of 10s*	0.57 [-2.80, 3.55]	-27.25 [-43.37, -21.06]	15.56 [14.59, 20.14]	<b>-21.46, 21.46</b>
30s*	-0.09 [-3.50, 3.32]	-20.71 [-26.61, -14.80]	20.52 [14.62, 26.43]	-21.46, 21.46
120s*	2.67 [0.86, 4.47]	-8.24 [-11.36, -5.12]	13.57 [10.45, 16.70]	-21.46, 21.46
<b>HF</b>				
120s	87.75 [8.48, 167.03]	-391.6 [-528.86, -254.25]	567.06 [429.76, 704.36]	<b>-318.86, 318.86</b>

Note. PR, mean pulse rate (bpm); HR, mean heart rate (bpm); RMSSD, root mean square of successive differences between normal IBIs (ms); SDNN, standard deviation of normal IBIs (ms); HF, high frequency (ms<sup>2</sup>); T = time interval; HR UST intervals distribution characteristics, all were homoscedastic and normally distributed with only the UST interval of 120s showing proportional bias; RMSSD UST distribution characteristics, all were homoscedastic with only the average of the three UST interval of 10s displaying non-normality which a log transformation could not alleviate and for the average of the three UST interval of 10s and the UST interval of 120s proportional bias was detected; SDNN UST intervals distribution characteristics, all were homoscedastic and normally distributed except for the UST interval of 10s which was heteroscedastic and the average of the three UST interval of 10s displaying non-normality which a log transformation could not alleviate; HF UST interval distribution characteristics, data was homoscedastic and normally distributed (see supplementary material E for the assumption handling); Median; 2.5 percentile, 97.5 percentile;  $\beta_0$ , the intercept;  $\beta_1$ , the slope coefficient; Bias SD, the SD of the residuals of the proportional bias model; GS; the gold standard value; LOA = 95% limits of agreement; 95% CI = 95% confidence interval; *A priori* 150% LOA, *a priori* 150% limits of agreement; Bold typeface = the 95% LOA is outside the *a priori* 150% LOA; \* indicates the proxies that were retained for the following step in the three-step hierarchical procedure.

Figure 1

Examples of Bland-Altman Plots Comparing Mean PR and PRV Metrics Obtained with UST Intervals with the E4 Wristband Versus Their 5min ECG Recording



Note. PR, mean pulse rate (bpm); HR, mean heart rate (i.e., in bpm); RMSSD, root mean square of successive differences between normal IBIs (i.e., in ms); SDNN, standard deviation of normal IBIs (i.e., in ms); HF, high frequency (i.e., in  $\text{ms}^2$ ); AVG10s, average of the three 10s intervals; density distribution of the differences at right side plot; 95% LOA, 95% limit of agreement; 95% CI, 95% confidence intervals; the *a priori* LOA of 150% only displayed if close to 95% LOA.

### 3. Step 3, magnitude of difference

Regarding mean PR, the effect sizes were negligible for all the UST intervals: T1 (Cliff's  $d = -.003$ , 95% CI [-.25, .25]), T2 (Cliff's  $d = .04$ , 95% CI [-.22, .28]), T3 (Cliff's  $d = .12$ , 95% CI [-.13, .37]), the average of three 10s intervals (Cliff's  $d = .05$ , 95% CI [-.21, .29]), 30s interval (Cliff's  $d = .001$ , 95% CI [-.25, .25]) and 120s interval (Cliff's  $d = .006$ , 95% CI [-.24, .26]). Concerning RMSSD, a small effect sizes was found for the UST interval of 30s (Cliff's  $d = -.20$ , 95% CI [-.43, .06]). Lastly, for the UST intervals estimating SDNN negligible effect sizes were found: the average of the three UST intervals of 10s (Cliff's  $d = .01$ , 95% CI [-.24, .26]), the UST interval of 30s (Cliff's  $d = -.03$ , 95% CI [-.28, .23]), and the UST interval of 120s (Cliff's  $d = -.12$ , 95% CI [-.36, .14]). These results show that the UST intervals specified above provide estimates of mean PR, RMSSD and SDNN comparable to the HR, RMSSD and SDNN values estimated with an ECG device 5min interval.

### E. Identifying and accommodating violations of assumptions Bland-Altman plot

*Assessment of violation of assumptions.* To assess if the **bias** can be approached as the mean of the observed differences, no proportional bias must be present. To measure the presence of proportional bias, we used the following equation (Bland & Altman, 1999):

$$(1) \text{bias}_i = \beta_0 + \beta_1 GS$$

Here  $\text{bias}_i$  represents the difference between the proxy and the gold standard.  $\beta_0$  is the intercept and  $\beta_1 GS$  represents the slope coefficient of the gold standard. If the slope coefficient is significant, proportional bias is detected, and the slope coefficient  $\beta_1 GS$  should represent the bias instead of the mean of the observed differences.

Next, to assess **scedasticity**, we applied the equation below (Bland & Altman, 1999). Namely, we regressed the gold standard on the absolute residuals (AR) of Eq. (2).

$$(2) AR_i = C_0 + C_1 GS_i$$

If the  $C_1 GS_i$  slope coefficient is significant, heteroscedasticity was detected, and homoscedasticity was assumed if the slope coefficient was insignificant.

Lastly, to measure the **normality** of the differences between the proxy and the gold standard, we used the Shapiro-Wilk test (Khatun, 2021; Menghini et al., 2021) and visually checked a Q-Q plot (Khatun, 2021).

*All assumptions met.* In the case of no proportional bias, homoscedasticity, and normality, we used the equation below:

$$(3) 95\% LOA = \text{bias} \pm 1.96SD$$

Here the bias depicts the mean of the observed differences between the proxy and the gold standard, and the 95% LOA is estimated by  $\pm 1.96SD$  from this bias. This is considered the standard procedure to represent the bias and its 95% LOA (Altman & Bland, 1983; Bland & Altman, 1999).

*Proportional bias.* If the bias was proportional and the data were homoscedastic and normally distributed, the computation of the 95% LOAs was as follows (Bland & Altman, 1999):

$$(4) 95\% LOA_i = \text{bias}_i \pm 1.96 \times \sqrt{\frac{1}{n} \sum_{i=1}^n [(x_{Pi} - x_{GSi}) - (\beta_0 + \beta_1 GS_i)]^2}$$

Here  $n$  is the sample size and  $x_{Pi} - x_{GSi}$  stands for the  $i$ -th proxy measure minus the  $i$ -th gold standard measure.

*Heteroscedasticity.* In the case of heteroscedasticity, and under the assumption of normality, the 95% LOAs were determined by the equation below (Bland & Altman, 1999).

$$(5) 95\% LOA_i = \text{bias} \pm 2.46 (c_0 + c_1 GS_i)$$

Here the type of bias used in the formula depends on whether it is proportional or not (i.e., mean as in Eq (3) or slope of Eq (1)).

*Non-normality.* If non-normality was detected, a log transformation was applied to the data (i.e., proxies and gold standard). If the non-normality was addressed, the mean represented the bias as in Eq. (3) or as a slope if the bias was proportional, as referred to in Eq. (1). If the log transformation did not alleviate the non-normality, we chose to represent the bias by the median of the observed differences instead of the mean. The median is a more suited central parameter as its less biased by the distribution of the observations than is the case for the mean as a central parameter (Leys et al., 2013). If the bias was proportional under these circumstances, we chose to use the equation below to represent the bias instead of Eq. (1):

$$(6) Q_{(\tau)} bias_i = \beta_{0(\tau)} + \beta_{1(\tau)} GS$$

Here,  $Q$  represents the conditional quantile estimation function and  $\tau$  is the quantile that needs to be estimated (i.e., range 0-1). To estimate the slope of the bias  $\tau = .50$  was used.  $\beta_{0(\tau)}$  and  $\beta_{1(\tau)}$  represent the intercept and slope coefficients. As such,  $\beta_{1(\tau)}$  represented the slope of the bias. The coefficients estimated by a quantile regression are less biased as they do not depend on the distribution of the observations, as is seen for standard linear regressions (Waldman, 2018).

If the non-normality was addressed, irrespective of the type of scedasticity, the 95% LOA was calculated following the equation below (Euser et al., 2008):

$$(7) 95\% LOA_i = bias \pm \{\forall x \in GS_i, antilog_{GS_i} = GS_i \times [(\pm 1) \times 2 \times \frac{e^{1.96SD(\ln bias_i) - 1}}{e^{1.96SD(\ln bias_i) + 1}}]\}$$

Here the type of bias used depended on whether it was proportional or not (i.e., mean as in Eq. (3) or slope of Eq. (1)).  $GS_i$  stands for the gold standard values and  $SD(\ln bias_i)$  represents the standard deviation of the differences between the log-transformed proxy and gold standard. So, for each element of the log-transformed gold standard values  $\forall x \in GS_i$ , the antilog slope  $antilog_{GS_i}$  of it is calculated.

If the non-normality was not alleviated, but no proportional bias was present, and the data were homoscedastic, we presented the 95% LOA as the 2.5 and 97.5 percentiles of the observed differences between the proxy and the gold standard (Bland & Altman, 1999). In case there was no proportional bias, but the data were heteroscedastic, we presented the 95% LOA following Eq. (7) relative to the bias (i.e., the median of the observed differences). In case non-normality was not alleviated with proportional bias and homoscedastic data, we used Eq. (6) to estimate the bias  $\beta_{1(\tau)}$  and Eq. (4) based on this proportional bias to obtain the 95% LOA. Under these same circumstances but with heteroscedastic data, we used Eq. (6) to estimate the bias and Eq. (7) to estimate the 95% LOA relative to this proportional bias.

## References

- Altman, D. G., & Bland, J. M. (1983). Measurement in medicine: The analysis of method comparison studies. *The Statistician*, 32(3), 307-317. doi: 10.2307/2987937
- Bland, J. M., & Altman, D. G. (1999). Measuring agreement in method comparison studies. *Statistical methods in medical research*, 8(2), 135–160. doi: 10.1177/096228029900800204
- Euser, A. M., Dekker, F. W., & le Cessie, S. (2008). A practical approach to Bland-Altman plots and variation coefficients for log transformed variables. *Journal of Clinical Epidemiology*, 61(10), 978-982. doi: 10.1016/j.jclinepi.2007.11.003

- Khatun, N. (2021). Applications of Normality Test in Statistical Analysis. *Open Journal of Statistics*, 11(01), 113. doi: 10.4236/ojs.2021.111006
- Leys, C., Ley, C., Klein, O., Bernard, P., & Licata, L. (2013). Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. *Journal of Experimental Social Psychology*, 49(4), 764–766. doi: 10.1016/j.jesp.2013.03.013
- Menghini, L., Cellini, N., Goldstone, A., Baker, F. C., & de Zambotti, M. (2021). Sleep-tracking technology: step-by-step guidelines and open-source code. *Sleep*, 44(2), 1-12. doi: 10.1093/sleep/zsaa170
- Waldmann, E. (2018). Quantile regression: A short story on how and why. *Statistical Modelling*, 18(3–4), 203–218. doi: 10.1177/1471082X18759142

## REFERENCES

- Ajayi, T. A., Salongo, L., Zang, Y., Wineinger, N., & Steinhubl, S. (2021). Mobile health-collected biophysical markers in children with serious illness-related pain. *Journal of Palliative Medicine*, 24(4), 580-588. doi: 10.1089/jpm.2020.0234
- Altman, D. G., & Bland, J. M. (1983). Measurement in medicine: The analysis of method comparison studies. *The Statistician*, 32(3), 307-317. doi: 10.2307/2987937
- Alqaraawi, A., Alwosheel, A., & Alasaad, A. (2016, May). Towards efficient heart rate variability estimation in artifact-induced Photoplethysmography signals. In *2016 IEEE Canadian Conference on Electrical and Computer Engineering (CCECE)* (pp. 1-6). IEEE. doi: 10.1109/CCECE.2016.7726853
- Association for the Advancement of Medical Instrumentation. (2002). Cardiac monitors, heart rate meters, and alarms. American National Standard (ANSI/AAMI EC13: 2002) Arlington, VA, 1-87
- Baek, H. J., Cho, C. H., Cho, J., & Woo, J. M. (2015). Reliability of ultra-short-term analysis as a surrogate of standard 5-min analysis of heart rate variability. *Telemedicine and e-Health*, 21(5), 404-414. doi: 10.1089/tmj.2014.0104
- Bent, B., Goldstein, B. A., Kibbe, W. A., & Dunn, J. P. (2020). Investigating sources of inaccuracy in wearable optical heart rate sensors. *Npj Digital Medicine*, 3(1), 1-9. doi.org/10.1038/s41746-020-0226-6
- Berntson, G. G., Thomas Bigger Jr, J., Eckberg, D. L., Grossman, P., Kaufmann, P. G., Malik, M., ... & Van Der Molen, M. W. (1997). Heart rate variability: origins, methods, and interpretive caveats. *Psychophysiology*, 34(6), 623-648. doi: 10.1111/j.1469-8986.1997.tb02140.x
- Bland, J. M., & Altman, D. G. (1999). Measuring agreement in method comparison studies. *Statistical methods in medical research*, 8(2), 135-160. doi: 10.1177/096228029900800204
- Campbell, J. I., & Thompson, V. A. (2012). MorePower 6.0 for ANOVA with relational confidence intervals and Bayesian analysis. *Behavior Research Methods*, 44(4), 1255-1265. doi: 10.3758/s13428-012-0186-0
- Castaneda, D., Esparza, A., Ghamari, M., Soltanpur, C., & Nazeran, H. (2018). A review on wearable photoplethysmography sensors and their potential future applications in health care. *International Journal of Biosensors & Bioelectronics*, 4(4), 195. doi: 10.15406/ijbsbe.2018.04.00125
- Charlot, K., Cornolo, J., Brugniaux, J. V., Richalet, J. P., & Pichon, A. (2009). Interchangeability between heart rate and photoplethysmography variabilities during sympathetic stimulations. *Physiological Measurement*, 30(12), 1357. doi: 10.1088/0967-3334/30/12/005
- Choi, K. H., Kim, J., Kwon, O. S., Kim, M. J., Ryu, Y. H., & Park, J. E. (2017). Is heart rate variability (HRV) an adequate tool for evaluating human emotions?—A focus on the use of the International Affective Picture System (IAPS). *Psychiatry Research*, 251, 192-196. doi: 10.1016/j.psychres.2017.02.025

- Euser, A. M., Dekker, F. W., & le Cessie, S. (2008). A practical approach to Bland-Altman plots and variation coefficients for log transformed variables. *Journal of Clinical Epidemiology*, *61*(10), 978-982. doi: 10.1016/j.jclinepi.2007.11.003
- Fallow, B. A., Tarumi, T., & Tanaka, H. (2013). Influence of skin type and wavelength on light wave reflectance. *Journal of Clinical Monitoring and Computing*, *27*(3), 313–317. doi.org/10.1007/s10877-013-9436-7
- Giavarina, D. (2015). Understanding Bland Altman analysis. *Biochemia Medica*, *25*(2), 141-151. doi: 10.11613/BM.2015.015
- Graham, S. A., Jeste, D. V., Lee, E. E., Wu, T. C., Tu, X., Kim, H. C., & Depp, C. A. (2019). Associations between heart rate variability measured with a wrist-worn sensor and older adults' physical function: observational study. *JMIR mHealth and uHealth*, *7*(10), e13757. doi: 10.2196/13757
- Heathers J. A. (2014). Everything Hertz: Methodological issues in short-term frequency-domain HRV. *Frontiers in physiology*, *5*, 177. doi: 10.3389/fphys.2014.00177
- Ishaque, S., Khan, N., & Krishnan, S. (2021). Trends in heart-rate variability signal analysis. *Frontiers in Digital Health*, *3*, 13. doi: 10.3389/fdgth.2021.639444
- Kiran Kumar, C., Manaswini, M., Maruty, K. N., Siva Kumar, A. V., & Mahesh Kumar, K. (2021). Association of heart rate variability measured by RR interval from ECG and pulse to pulse interval from photoplethysmography. *Clinical Epidemiology and Global Health*, *10*(100698), 2213-3984. doi: 10.1016/j.cegh.2021.100698
- Kuipers, M., Richter, M., Scheepers, D., Immink, M. A., Sjak-Shie, E., & van Steenbergen, H. (2017). How effortful is cognitive control? Insights from a novel method measuring single-trial evoked beta-adrenergic cardiac reactivity. *International Journal of Psychophysiology*, *119*, 87-92. doi: 10.1016/j.ijpsycho.2016.10.007
- Krouwer, J. S. (2008). Why Bland-Altman plots should use X, not  $(Y + X)/2$  when X is a reference method. *Statistics in Medicine*, *27*(5), 778–780. doi.org/10.1002/sim.3086
- Kumral, D., Schaare, H. L., Beyer, F., Reinelt, J., Uhlig, M., Liem, F., ... Gaebler, M. (2019). The age-dependent relationship between resting heart rate variability and function brain connectivity. *NeuroImage*, *185*, 521-533. doi: 10.1016/j.neuroimage.2018.10.027
- Laborde, S., Mosley, E., & Thayer, J. F. (2017). Heart rate variability and cardiac vagal tone in psychophysiological research—recommendations for experiment planning, data analysis, and data reporting. *Frontiers in Psychology*, *8*, 213. doi: 10.3389/fpsyg.2017.00213
- Lackner, H. K., Weiss, E. M., Schuler, G., Hinghofer-Szalkay, H., Samson, A. C., & Papousek, I. (2013). I got it! Transient cardiovascular response to the perception of humor. *Biological Psychology*, *93*(1), 33-40. doi: 10.1016/j.biopsycho.2013.01.014
- Lipponen, J. A., & Tarvainen, M. P. (2019). A robust algorithm for heart rate variability time series artefact correction using novel beat classification. *Journal of Medical Engineering & Technology*, *43*(3), 173-181. doi: 10.1080/03091902.2019.1640306
- Lu, G., & Yang, F. (2009). Limitations of oximetry to measure heart rate variability measures. *Cardiovascular Engineering*, *9*(3), 119-125. doi: 10.1007/s10558-009-9082-3



- Lu, S., Zhao, H., Ju, K., Shin, K., Lee, M., Shelley, K., & Chon, K. H. (2008). Can photoplethysmography variability serve as an alternative approach to obtain heart rate variability information? *Journal of Clinical Monitoring and Computing*, 22(1), 23–29. doi: 10.1007/s10877-007-9103-y
- McCarthy, C., Pradhan, N., Redpath, C., & Adler, A. (2016, May). Validation of the Empatica E4 wristband. In *2016 IEEE EMBS International Student Conference (ISC)* (pp. 1-4). IEEE. doi: 10.1109/EMBSISC.2016.7508621
- Mejía-Mejía, E., May, J. M., Torres, R., & Kyriacou, P. A. (2020). Pulse rate variability in cardiovascular health: A review on its applications and relationship with heart rate variability. *Physiological Measurement*, 41(7), 07TR01. doi: 10.1088/1361-6579/ab998c
- Menghini, L., Cellini, N., Goldstone, A., Baker, F. C., & de Zambotti, M. (2021). Sleep-tracking technology: step-by-step guidelines and open-source code. *Sleep*, 44(2), 1-12. doi: 10.1093/sleep/zsaa170
- Menghini, L., Gianfranchi, E., Cellini, N., Patron, E., Tagliabue, M., & Sarlo, M. (2019). Stressing the accuracy: Wrist-worn wearable sensor validation over different conditions. *Psychophysiology*, 56(11), e13441. doi: 10.1111/psyp.13441
- Milstein, N., & Gordon, I. (2020). Validating measures of electrodermal activity and heart rate variability derived from the Empatica E4 utilized in research settings that involve interactive dyadic States. *Frontiers in behavioral neuroscience*, 14, 148. doi: 10.3389/fnbeh.2020.00148
- Mishra, T., Wang, M., Metwally, A. A., Bogu, G. K., Brooks, A. W., Bahmani, A., ... & Snyder, M. P. (2020). Pre-symptomatic detection of COVID-19 from smartwatch data. *Nature biomedical engineering*, 4(12), 1208-1220. doi: 10.1038/s41551-020-00640-6
- Mol, A., Slangen, L. R. N., Trappenburg, M. C., Reijnerse, E. M., van Wezel, R. J. A., Meskers, C. G. M., & Maier, A. B. (2020). Blood pressure drop rate after standing up is associated with frailty and number of falls in geriatric outpatients. *Journal of the American Heart Association*, 9(7), e014688. doi: 10.1161/JAHA.119.014688
- Munoz, M. L., van Roon, A., Riese, H., Thio, C., Oostenbroek, E., Westrik, I., ... & Snieder, H. (2015). Validity of (ultra-) short recordings for heart rate variability measurements. *PloS One*, 10(9), e0138921. doi: 10.1371/journal.pone.0138921
- Nuske, H. J., Goodwin, M. S., Kushleyeva, Y., Forsyth, D., Pennington, J. W., Masino, A. J., ... & Herrington, J. D. (2021). Evaluating commercially available wireless cardiovascular monitors for measuring and transmitting real-time physiological responses in children with autism. *Autism Research*, 15(1), 117-130. doi: 10.1002/aur.2633
- Nussinovitch, U., Elishkevitz, K. P., Katz, K., Nussinovitch, M., Segev, S., Volovitz, B., & Nussinovitch, N. (2011). Reliability of ultra-short ECG indices for heart rate variability. *Annals of Noninvasive Electrocardiology*, 16(2), 117-122. doi: 10.1111/j.1542-474X.2011.00417.x
- Ollander, S., Godin, C., Campagne, A., & Charbonnier, S. (2016, October). A comparison of wearable and stationary sensors for stress detection. In *2016 IEEE International Conference on Systems, Man, and Cybernetics (SMC)* (pp. 004362-004366). IEEE. doi: 10.1109/SMC.2016.7844917

- Ottaviani, C., Zingaretti, P., Petta, A. M., Antonucci, G., Thayer, J. F., & Spitoni, G. F. (2018). Resting heart rate variability predicts inhibitory control above and beyond impulsivity. *Journal of Psychophysiology*, 33. doi: 10.1027/0269-8803/a000222
- Pecchia, L., Castaldo, R., Montesinos, L., & Melillo, P. (2018). Are ultra-short heart rate variability features good surrogates of short-term ones? State-of-the-art review and recommendations. *Healthcare Technology Letters*, 5(3), 94-100. doi: 10.1049/htl.2017.0090
- Pichon, A., Roulaud, M., Antoine-Jonville, S., de Bisschop, C., & Denjean, A. (2006). Spectral analysis of heart rate variability: interchangeability between autoregressive analysis and fast Fourier transform. *Journal of electrocardiology*, 39(1), 31-37. doi: 10.1016/j.jelectrocard.2005.08.001
- Pulopulos, M. M., Vanderhasselt, M. A., & De Raedt, R. (2018). Association between changes in heart rate variability during the anticipation of a stressful situation and the stress-induced cortisol response. *Psychoneuroendocrinology*, 94, 63-71. doi: 10.1016/j.psyneuen.2018.05.004
- Quintana, D. S., & Heathers, J. A. J. (2014). Considerations in the assessment of heart rate variability in biobehavioral research. *Frontiers in Psychology*, 5(JUL), 1–10. doi: doi.org/10.3389/fpsyg.2014.00805
- R Core Team. (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. url: R-project.org
- Rahman, J. S., Gedeon, T., Caldwell, S., Jones, R., & Jin, Z. (2021). Towards effective music therapy for mental health care using machine learning tools: human affective reasoning and music genres. *Journal of Artificial Intelligence and Soft Computing Research*, 11, 5-20. doi: 10.2478/jaiscr-2021-0001
- Romano, J., Kromrey, J. D., Coraggio, J., Skowronek, J., & Devine, L. (2006, October). Exploring methods for evaluating group differences on the NSSE and other surveys: Are the t-test and Cohen'sd indices the most appropriate choices. In *annual meeting of the Southern Association for Institutional Research* (pp. 1-51).
- Ryan, W., Conigrave, J. H., Basarkod, G., Ciarrochi, J., & Sahdra, B. K. (2019). When is it good to use wristband devices to measure HRV?: Introducing a new method for evaluating the quality of data from photoplethysmography-based HRV devices. doi: 10.31234/osf.io/t3gdz
- Salahuddin, L., Cho, J., Jeong, M. G., & Kim, D. (2007, August). Ultra short term analysis of heart rate variability for monitoring mental stress in mobile settings. In *2007 29th annual international conference of the IEEE Engineering in Medicine and Biology Society* (pp. 4656-4659). IEEE. doi: 10.1109/IEMBS.2007.4353378
- Schuermans, A. A., de Loeff, P., Nijhof, K. S., Rosada, C., Scholte, R. H., Popma, A., & Otten, R. (2020). Validity of the Empatica E4 wristband to measure heart rate variability (HRV) parameters: A comparison to electrocardiography (ECG). *Journal of Medical Systems*, 44(11), 1-11. doi: 10.1007/s10916-020-01648-w
- Shaffer, F., & Ginsberg, J. P. (2017). An overview of heart rate variability metrics and norms. *Frontiers in Public Health*, 5, 258. doi: 10.3389/fpubh.2017.00258

- Shaffer, F., McCraty, R., & Zerr, C. L. (2014). A healthy heart is not a metronome: an integrative review of the heart's anatomy and heart rate variability. *Frontiers in Psychology, 5*, 1040. doi: 10.3389/fpsyg.2014.01040
- Shaffer, F., Meehan, Z. M., & Zerr, C. L. (2020). A critical review of ultra-short-term heart rate variability norms research. *Frontiers in neuroscience, 14*, 594880. doi: 10.3389/fnins.2020.594880
- Shaffer, F., Shearman, S., & Meehan, Z. M. (2016). The promise of ultra-short-term (UST) heart rate variability measurements. *Biofeedback, 44*(4), 229-233. doi: 10.5298/1081-5937-44.3.09
- Shen, W., Yuan, Y., Tang, C., Shi, C., Liu, C., Luo, J., & Zhang, X. (2017). In search of somatic precursors of spontaneous insight. *Journal of Psychophysiology, 32*(3), 97–105. doi: 10.1027/0269-8803/a000188
- Silverthorn, D.U. (2004). *Human Physiology: An Integrated Approach with Interactive Physiology, Third Edition*. San Francisco: Pearson Education, Inc.
- So, T. Y., Li, M. Y. E., & Lau, H. (2021). Between-subject correlation of heart rate variability predicts movie preferences. *PloS One, 16*(2), e0247625. doi: 10.1371/journal.pone.0247625
- Stone, J. D., Ulman, H. K., Tran, K., Thompson, A. G., Halter, M. D., Ramadan, J. H., ... & Hagen, J. A. (2021). Assessing the accuracy of popular commercial technologies that measure resting heart rate and heart rate variability. *Frontiers in Sports and Active Living, 3*, 37. doi: 10.3389/fspor.2021.585870
- Tarvainen, M. P., Lipponen, J., Niskanen, J. P., & Ranta-aho, P. O. (2020). Kubios HRV User Guide, Kubios Oy.
- Tarvainen, M. P., Niskanen, J. P., Lipponen, J. A., Ranta-Aho, P. O., & Karjalainen, P. A. (2014). Kubios HRV–heart rate variability analysis software. *Computer Methods and Programs in Biomedicine, 113*(1), 210-220. doi: 10.1016/j.cmpb.2013.07.024
- Tarvainen, M. P., Ranta-Aho, P. O., & Karjalainen, P. A. (2002). An advanced detrending method with application to HRV analysis. *IEEE Transactions on Biomedical Engineering, 49*(2), 172-175. doi: 10.1109/10.979357
- Task Force of the European Society of Cardiology the North American Society of Pacing Electrophysiology (1996). Heart rate variability: Standards of measurement, physiological interpretation, and clinical use. *Circulation, 93*(5), 1043–1065. doi: 10.1161/01.CIR.93.5.1043
- Torchiano, M. (2020). *effsize: Efficient Effect Size Computation* (R package version 0.8.1). doi: 10.5281/zenodo.1480624
- Umetani, K., Singer, D. H., McCraty, R., & Atkinson, M. (1998). Twenty-four hour time domain heart rate variability and heart rate: Relations to age and gender over nine decades. *Journals of the American College of Cardiology, 31*(3), 593-601. doi: 10.1016/s0735-1097(97)00554-8

- van Lier, H. G., Pieterse, M. E., Garde, A., Postel, M. G., de Haan, H. A., Vollenbroek-Hutten, M. M., ... & Noordzij, M. L. (2020). A standardized validity assessment protocol for physiological signals from wearable technology: Methodological underpinnings and an application to the E4 biosensor. *Behavior Research Methods*, 52(2), 607–629. doi: 10.3758/s13428-019-01263-9
- Yuda, E., Shibata, M., Ogata, Y., Ueda, N., Yambe, T., Yoshizawa, M., & Hayano, J. (2020). Pulse rate variability: A new biomarker, not a surrogate for heart rate variability. *Journal of Physiological Anthropology*, 39(1), 21. doi: 10.1186/s40101-020-00233-x

#### **AUTHOR NOTE**

#### **ACKNOWLEDGMENTS**

We thank Febe Demeyer, Christo Bratanov, Yujing Liang, and Amar Music for their assistance with the data collection and their critical thought and the Fonds de la Recherche Scientifique and the Research Foundation Flanders (FWO) for providing the opportunity to conduct this research under a research fellow grant.

#### **FUNDING**

This work was supported by the "Fonds de la Recherche Scientifique" [grant number 34736358, 2019] and the Research Foundation Flanders [grant number G096919N].