

1 Article

2 About the Accuracy and Problems of Consumer 3 Devices in the Assessment of Sleep

4 Mohamed S. Ameen ¹, Lok Man Cheung ^{1,2}, Theresa Hauser ¹, Michael Hahn ¹, Manuel Schabus ^{1,3,*}

5 ¹ University of Salzburg, Department of Psychology, Laboratory for Sleep, Cognition and Consciousness
6 Research, Hellbrunner Strasse 34, 5020 Salzburg.

7 ² School of Psychology, University of Surrey, Surrey, United Kingdom, Stag Hill, Guildford GU2 7XH, UK.

8 ³ Center for Cognitive Neuroscience Salzburg (CCNS), Hellbrunner Strasse 34, 5020 Salzburg.

9 * Correspondence: manuel.schabus@sbg.ac.at; Tel: (+43) 66280445113

10 **Abstract:** Commercial sleep devices and mobile-phone applications for scoring sleep are gaining
11 ground. In order to provide reliable information about the quantity and/or quality of sleep, their
12 performance needs to be assessed against the current gold-standard, i.e. polysomnography (PSG;
13 measuring brain, eye and muscle activity). We here assessed some commercially available sleep
14 trackers, namely; a commercial activity tracker: Mi band (Xiaomi, BJ, CHN), a scientific actigraph:
15 Motionwatch 8 (CamNTEch, CB, UK), and a much used sleep application: Sleep Cycle (Northcube,
16 GOT, SE). We recorded 27 nights in healthy sleepers using PSG and these devices. Surprisingly, all
17 devices had very poor agreement with the gold standard. Sleep parameter comparisons revealed
18 that specifically the Mi band and the sleep cycle application had difficulties in detecting wake
19 periods which negatively affected the total sleep time and sleep efficiency estimations. However, all
20 3 devices were good in detecting the most basic parameter, the actual time in bed. In summary, our
21 results suggest that, to-date; available sleep trackers do not provide meaningful sleep analysis but
22 may be interesting for simply tracking times in bed. A much closer interaction with the scientific
23 field seems necessary if reliable information shall be derived from such devices in the future.

24 **Keywords:** Wrist-worn devices; Sleep trackers; Activity trackers; Sleep classification;
25 Polysomnography.

26

27 1. Introduction

28 Our knowledge about the structure and function of sleep is derived mainly from recordings that
29 are done in sleep laboratories. In these recordings physiological activity is measured using
30 polysomnography (PSG), which requires a combination of electroencephalography (EEG),
31 electrooculography (EOG) and electromyography (EMG) data. Although these recordings
32 contribute widely to our constantly expanding knowledge about sleep, their major drawback is that
33 they do not mimic the habitual sleeping environment at home. The effect of the laboratory setup on
34 sleep has been repeatedly addressed and several studies have highlighted differences in sleep
35 parameters between at-home and laboratory sleep recordings. For instance, Portier and colleagues
36 [1] compared several sleep parameters between a night of sleep in the laboratory versus a night of
37 sleep at-home. Specifically, the authors reported a significant reduction in the amount of total sleep

38 time (TST) and time in bed (TiB) as well as the deterioration of the subjective experience of sleepers
39 on the quality of their sleep in the laboratory as compared to at-home sleep. In a similar vein,
40 home-based PSG recording showed higher sleep efficiency (SE) values than hospital-based PSG for
41 the same participants [2]. These findings suggest that in-laboratory sleep does not accurately reflect
42 habitual sleep (at-home sleep) and consequently might introduce some bias in the diagnosis of some
43 sleep disorders. Indeed, it has been shown that home-based PSG recordings show better dissociation
44 between healthy sleepers and insomnia sufferers than in laboratory settings [3]. Therefore, it is a
45 priority in the field of sleep research and sleep medicine to develop better tools that can accurately
46 and reliably measure sleep at home and in a vast amount of the general population.

47 Already in the last years we have witnessed a vast increase in the available consumer devices, i.e.
48 sleep devices and mobile phone applications, which aim to assess and ultimately improve sleep.
49 These devices might be of potential help to overcome the bias induced by the laboratory setting as
50 they are supposed to assess sleep outside the laboratory with minimal effort for the end-user.
51 However, it is essential to scientifically test and compare these devices against the “gold-standard”
52 to ensure that such devices and applications do not provide random feedback to the naïve end-user
53 and potentially backfire with “unintended effects on sleep beliefs and behaviors” [4]. Only the
54 adherence of such devices and applications to the gold standard ensures reliability and validity and
55 ethically justifies these new methods advertised by the industry.

56 The aim of this study was therefore to assess the performance of some of these readily used
57 consumer devices which claim to monitor sleep and provide reliable information about sleep quality
58 and sleep architecture night-by-night. Specifically we assessed sleep data from 2 devices: 1) a
59 commercial activity tracker, the Mi band (v2, Xiaomi, BJ, CHN), 2) a scientific actigraph watch,
60 Motionwatch 8 (CamNTech, CB, UK) as well as one readily used mobile phone application: the Sleep
61 Cycle App (v3.0.1.2511-release; Northcube, GOT, SE). We compared the full set of sleep measures
62 from these 3 platforms against our PSG gold standard and relied on semi-automatic sleep staging
63 using the SOMNOlyzer 24X7 solution [5,6].

64 2. Materials and Methods

65 Study sample: For the study we recruited 19 healthy participants (13 females, mean age: 29±13.
66 Range: 19-64 years). Participants arrived at the sleep laboratory of the University of Salzburg at 9pm.
67 They were instructed about the procedure and the purpose of the experiment. After signing the
68 consent forms they were given the Mi band (MB) and the MotionWatch (MW). After we confirmed
69 that both the MB and the MW were recording we started with the PSG preparation. Before turning
70 lights off and starting the PSG recording we started the Sleep Cycle application (SC) and placed the
71 device next to the subject. Participants went to bed at around 11 pm and stayed in bed (TiB, time in
72 bed) for approximately 8h (452.29 ± 81.78 min). Two of these participants had to be excluded from
73 our analyses due to technical problems with the PSG recordings. Additionally, we recorded 8
74 ambulatory “home” PSG nights using an ambulatory EEG device together with the MB device and
75 the SC application after the participants visited the lab for electrode placement; therefore, we had a
76 total of 27 nights of PSG recordings. Due to other problems with some of the devices and
77 applications we finally analyzed 21 nights for the MB, and 12 nights for the MW and the SC.

78 EEG Data Acquisition: For the nights spent in the laboratory, brain activity was recorded using
79 high-density-EEG with a 256-electrode GSN HydroCel Geodesic Sensor Net (Electrical 478
80 Geodesics Inc., Eugene, Oregon, USA) and a Net Amps 400 amplifier. Additionally, we recorded
81 electrocardiography (ECG), electromyography (EMG) and electrooculography (EOG) using bipolar
82 electrodes. Ambulatory PSG was recorded using a 16-channel EEG, bipolar EMG and EOG using the
83 AlphaEEG amplifier and NeuroSpeed software (Alpha Trace Medical Systems, Vienna, Austria).

84 Sleep Scoring: Our PSG was analyzed for sleep stages using the computer-assisted sleep
85 classification system Somnolyer 24x7 as developed by the SIESTA group (The SIESTA Group
86 Schlafanalyse GmbH, Vienna, Austria; [5,6]) and was following the revised standard criteria
87 described by the American Association for Sleep Medicine (AASM, [7]). The derived sleep features
88 and sleep stages serve as gold-standard for the rest of the analyses. Sleep staging for the SC
89 application was realized via a simple image processing of the figures generated by the application;
90 basically we discretized the SC illustrations into 3 sleep-wake states as suggested by the application
91 in wake, light sleep and deep sleep (cf. Suppl. Material and Supplementary Figure S1 for more
92 details).

93 Statistical analysis: The following five main sleep parameters were evaluated: i) sleep onset
94 latency (SOL), ii) sleep efficiency (SE), iii) wake after sleep onset (WASO), iv) total sleep time (TST),
95 and v) time in bed (TiB). SOL was defined as the difference between the start of the recording and
96 the time when the participant actually fell asleep (i.e. the 1st N1 or “light sleep” epoch). SE (%) was
97 defined as: $(TST / TiB) * 100$. Importantly, measurements from all the devices were accurately
98 synchronized to the start of the PSG recording. Correlations were computed non-parametrically
99 using spearman correlations. For the MB device measurements we needed to calculate SOL, SE, and
100 TST (as described above) manually and used WASO values as provided by the device. For the SC,
101 we manually calculated SOL, SE, TST, and WASO (as described above) and used the time points
102 provided by the application to calculate TiB. For the MW measurements, the sensitivity threshold
103 was set to 20 activity counts and adjusted the lights off and lights on times according to sleep diaries
104 as usually done for scientific actigraph measurements. SOL, SE, TST, WASO and TiB values are then
105 used as provided by the MotionWare software (v1.1.20, empire Software GmbH, Cologne,
106 Germany).

107 Bland-Altman plots were used to quantify the agreement between the PSG gold standard and
108 the three consumer devices. The measured bias is defined as the mean of the difference between the
109 two-paired measurements. That is, the further this value is from zero, i.e. the line of equality
110 (difference = 0), the higher the error in the measurement. Spearman-correlations are used to illustrate
111 systematic linear biases of the devices, and are reported at $p < .01$ (corrected for the 5 dependent sleep
112 variables analyzed).

113 Epoch-wise comparison of sleep stages: For analyzing the epoch-by-epoch agreement of the
114 gold-standard with the three consumer devices we always synchronized the recording start with the
115 start of the PSG recording. In case one device started recording after the other (for example, PSG
116 after SC or vice versa) we simply discarded the earlier epochs and started the analysis from the first
117 epoch which was scored by both. As mentioned above we used the graphs provided by the MB

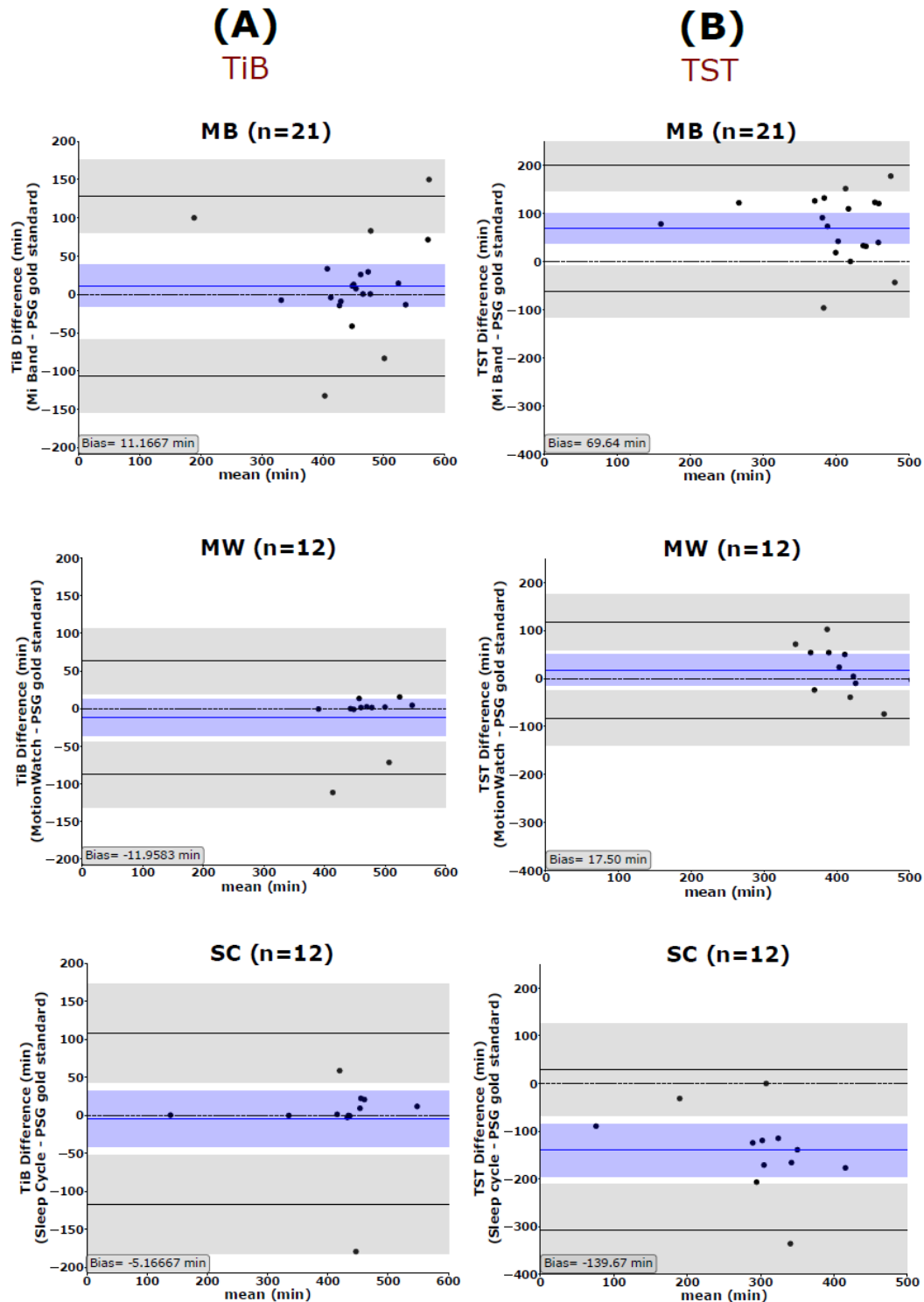
118 device and the SC application in order to divide sleep/waking into awake, light and deep sleep and
119 in 30s epochs. For the PSG gold standard light sleep was defined as stages N1 and N2 while Deep
120 sleep is defined as N3 stages. Importantly, we excluded PSG epochs which were scored as stage
121 REM according to the AASM from the analysis as all 3 devices and applications provide no
122 information about REM (or “dreaming”) sleep. We report two main parameters for the epoch-wise
123 agreement; sensitivity and positive predictive value (PPV). Sensitivity (in %) estimates the
124 epoch-by-epoch agreement between the MB and SC with the gold standard by measuring the % of
125 correct classifications (according to the PSG standard) per sleep stage (that is, for example labelling
126 79% of all light sleep detections by the PSG as “light sleep”). The positive predictive value (PPV), on
127 the other hand, is the probability that the assigned state (by the device or app) is indeed that specific
128 state in the gold standard (that is, for example only 41% of assigned “light sleep” epochs are actually
129 light sleep epochs and no other sleep states). Cohen’s Kappa (K) was used to assess the pairwise
130 agreement between the devices. Epoch-wise analysis was computed in SPSS software (IBM Corp.
131 Released 2017. IBM SPSS Statistics for Windows, Version 25.0. Armonk, NY: IBM Corp.). Cohen’s
132 Kappa scores between 0.21–0.40 is often considered as fair, 0.41–0.60 as moderate, and 0.61–0.80 as
133 substantial agreement according to Landis and Koch (1977, [7]).

134 3. Results

135 The mean values of the key features of sleep across all participants according to the PSG gold
136 standard were 434.58 ± 95.83 minutes for the TiB, 370.12 ± 104.43 minutes for the TST, $84.08 \pm$
137 13.22% for the SE, 25.98 ± 19.35 minutes for the SOL and 39.08 ± 38.43 minutes for WASO. As a first
138 analysis we simply checked whether the mean sleep values per participant and night correlate
139 between the gold standard and the devices. For TiB we found good agreement, that is significant
140 positive associations, of the gold standard values with the 3 consumer devices (MB: $r= 0.72$,
141 $p=0.0002$; SC: $r= 0.67$, $p=0.02$; MW: $r= 0.77$, $p=0.03$). For TST we only found one moderately positive
142 association for the MB device ($r= 0.49$, $p= 0.02$), while MW was the only device that showed a
143 significant positive correlation for WASO time ($r= 0.78$, $p= 0.02$) (see supplementary figures S2-5).
144 This low agreement is already surprising given that these are the simple associations of the mean
145 values per subject, e.g., whether people taking longer to fall asleep in the case of SOL measurements
146 (according to the PSG gold standard) also tend to fall asleep later according to the output of one of
147 the consumer devices.

148 3.1. Bland Altman plots:

149 We used the Bland and Altman analysis to visualize the degree of agreement between the PSG gold
150 standard and each of the 3 aforementioned devices/applications (cf. Figures 1-3). The most global
151 key features of sleep, namely TIB and TST are depicted in Figure 1. Looking at the mean difference,
152 there is only a slight bias towards over- or underestimating TIB (MB: 11.1 ± 59.96 min, MW: $11.96 \pm$
153 38.48 min & SC: -5.17 ± 57.57 min). However, the 95% confidence interval also indicates that for
154 single cases the devices may still over- or underestimate TIB by an hour or more (cf. Figure 1.A).
155 Mean TST is systematically overestimated by the MB by more than an hour (69.64 ± 67.43 min), and
156 underestimated by the SC application by more than two hours on average (-139.67 ± 85.87 min)
157 (Figure 1.B).



158

159

160

161

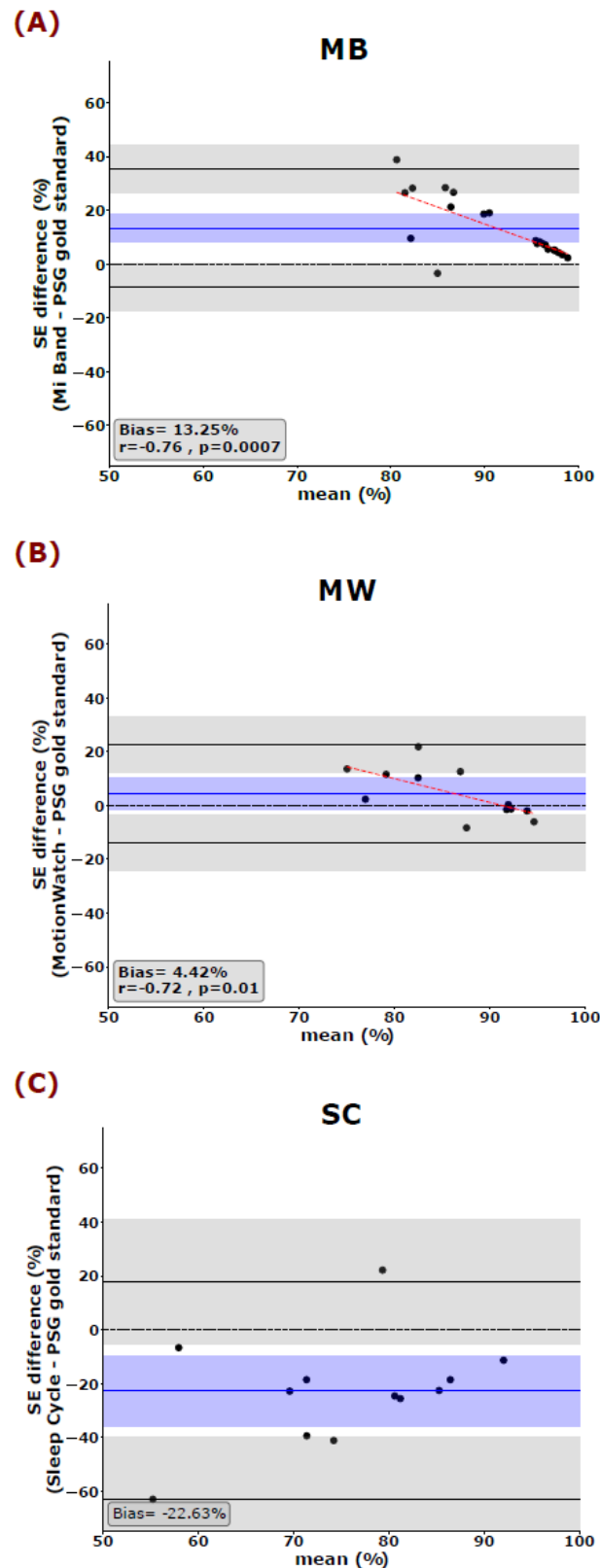
162

163

164

Figure 1. Bland Altman plots show the agreement of the MB, the MW and the SC with the PSG in measuring (A) TiB but not (B) TST. The blue horizontal line represents the mean difference between the two measurements and the shaded blue area represents the 95% CI of the mean difference. Black horizontal lines mark the 1.96SD from the mean and the grey shadings represent their 95%CI. The black dashed line is the line of equality (difference=0). TiB: time in bed, TST: total sleep time, MB: Mi Band, MW: MotionWatch and SC: sleep cycle application.

165 All 3 tested devices/applications showed inaccuracies in estimating SE, with the scientific MW
166 giving the best results and no systematic over- or underestimation of SE (Figure 2). The mean
167 differences indicate that the MB systematically overestimated SE (13.25%) whereas the SC
168 application systematically underestimated (22.63%) SE. Interestingly, spearman correlations
169 indicated that the MB linearly shows greater errors the worse the real PSG gold standard SE was,
170 that is the MB has a strong bias towards quantifying SE better than it is.



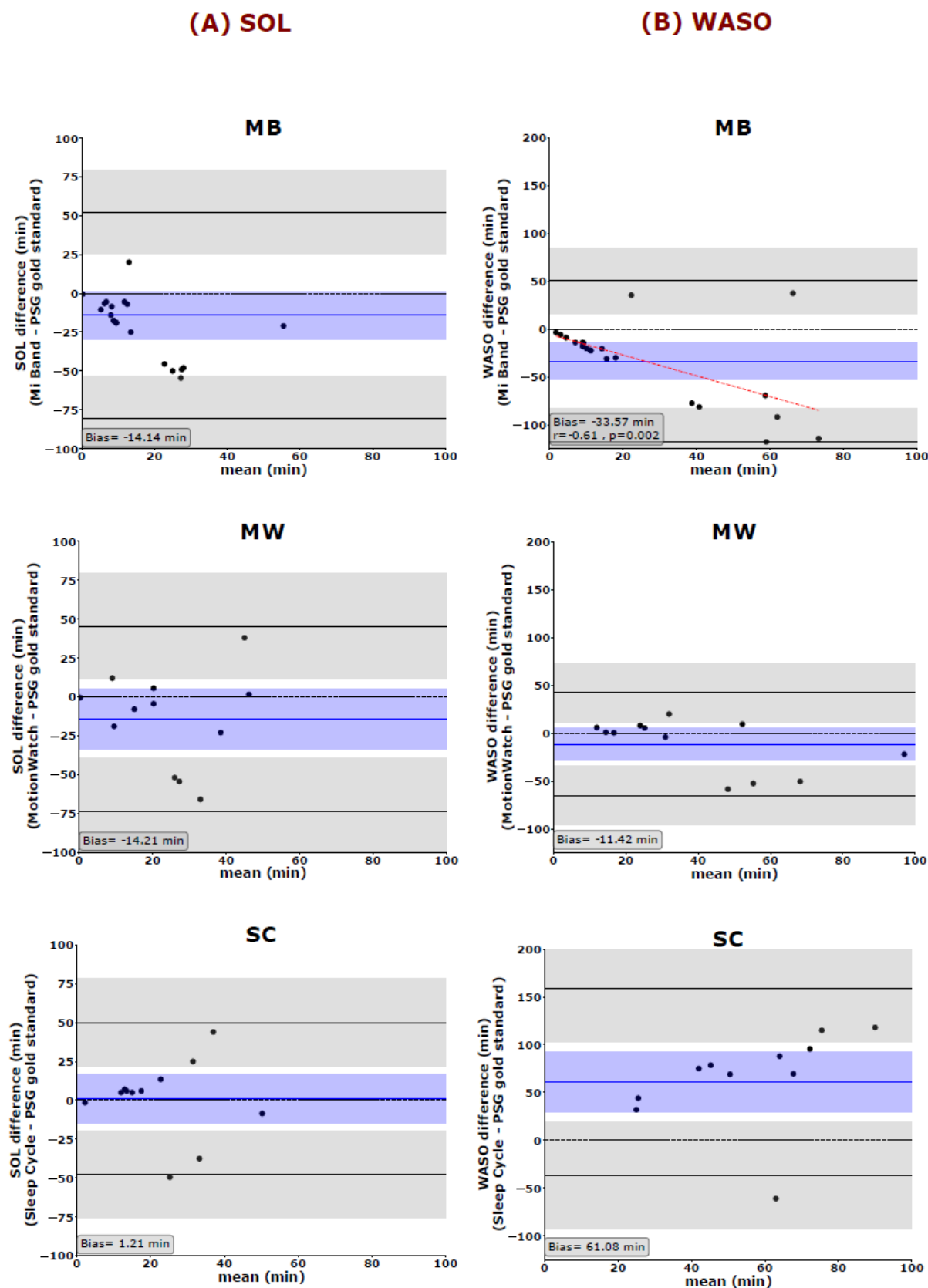
171 **Figure 2. Bland Altman plots show that the PSG gold standard with the (A) MB, the (B) MW**
172 **and (C) the SC in measuring sleep efficiency (SE).** The blue horizontal line represents the mean
173 difference between the two measurements and the shaded blue area represents the 95%-CI of the
174 mean difference. Black horizontal lines mark the 1.96SD from the mean and the grey shadings
175 represent their 95%CI. The black dashed line is the line of equality (difference=0) and the red
176 dashed line represents the spearman-correlation between the difference and the average of the
177 two measurements. SE: Sleep Efficiency, MB: mi band, MW: MotionWatch and SC: sleep cycle
178

179 Moreover, we observed a systematic error in the estimation of the WASO time by the MB and the SC
180 but not the MW (Figure 3.A). While the MB device underestimates WASO (-33.57 ± 42.84 min) the SC
181 App overestimates WASO systematically (61.10 ± 49.90 min). In addition there is a linear trend in the
182 data showing that the MB underestimates WASO time the more the longer actual WASO time gets.
183 The mean difference of the 3 devices/applications to the gold standard is closer to zero for SOL,
184 however it is to be noted that here also the range of possible values is much more limited ($36.18 \pm$
185 38.37 min for the gold standard). Only the MB shows a linear trend with stronger underestimation of
186 SOL the longer SOL actually was (in the gold standard) (Figure 3.B).

187 3.1. Epoch-wise agreement per sleep stage

188 Table 1 shows the overall and stage-wise agreement between the 30s-epochs scored by our PSG gold
189 standard and both the MB device and the SC application. Note that the MW is disregarded in this
190 respect as standard MW outputs do not provide (or claim to allow) sleep staging classifications. The
191 overall agreement over all epochs from all subjects (16350 epochs for the MB, 11243 epochs for the
192 MW and 9504 epochs for the SC) between the gold standard PSG scoring and the MB was relatively
193 low (53.31% , $k=0.14$) and even lower for the SC device (46.34%, $k=0.18$). Table 1 also illustrates that
194 the highest level of agreement for the MB was in determining light sleep (sensitivity = 70.6% and
195 PPV = 57.8%) and the lowest sensitivity for the MB was for detecting wakefulness (sensitivity = 5.5%;
196 PPV= 62.8%). Conversely, however, SC had moderate sensitivity in identifying awake epochs
197 (sensitivity= 55.6%) and an unacceptable PPV value of 24.3%, meaning that only 24.3% of wake
198 classified epochs are indeed wake according to the PSG gold-standard. On the other hand, SC had
199 low sensitivity in detecting light sleep (40.9%) yet when it classified light sleep this was the true state
200 in 61.2% of the cases (i.e., PPV= 61.2%). Moreover, for “deep sleep” classification we found very poor
201 performance for the MB (sensitivity= 47.2%, PPV= 43.6%) and poor performance for the SC App
202 (sensitivity= 52.0%, PPV= 53.0%).

203 Given this poor performance in correctly classifying sleep stages we then investigated the ability
204 of these devices and the SC application to simply differentiate between sleep (light sleep, deep sleep
205 or REM) and wakefulness and included the scientific MW device (whose software anyways only
206 provides wake and sleep classes). Here we then found good overall agreement for the MB and the
207 MW (>80%, cf. Table 2), and rather poor overall agreement for the SC App (65.9%). Kappa pairwise
208 agreement indicates a “fair” agreement for the MW but poor agreements for the other two
209 devices/apps. Specifically, the output shows that the MB and MW devices on the arm wrist are very
210 good when only “sleep” detection is needed (MB: sensitivity= 99.5%, PPV= 86.8%; MW: sensitivity=
211 92.9%, PPV= 88.2%). Severe difficulties remain in assigning stage wake in all three devices/apps and
212 therefore a proper estimation of overall sleep efficiency or sleep quality remains non feasible.



213

214

215

216

Figure 3. Bland-Altman plots of the SOL and WASO measurements showing differences between the PSG gold standard and the MB, the MW and the SC. The blue horizontal line represents the mean difference between the two measurements and the shaded blue area represents the 95%-CI of the

217 mean difference. Black horizontal lines mark the 1.96SD from the mean and the grey shadings
 218 represent their 95% CI. The black dashed line is the line of equality (difference=0) and the red dashed
 219 line represents the spearman-correlation between the difference and the average of the two
 220 measurements. SOL: Sleep onset latency, WASO: wake after sleep onset. MB: Mi band, MW:
 221 MotionWatch and SC: sleep cycle application.

222
 223 Although the SC application is as good as the wrist band devices in assigning “sleep” to an
 224 epoch, that is then the app is correct in 91.3% of these cases, it still misses a third of all sleep
 225 epochs (sensitivity= 67.4%). Importantly, however, in the case of the MB and the SC, OA
 226 increases while Kappa scores dropped when we pooled all sleep stages in one stage which
 227 likely indicates serious bias in the scoring algorithms of the MB device and the SC app.

228 **Table 1.** Percentages of agreement table for 3 stages sleep scoring (Awake/ Light sleep/ Deep sleep) between the
 229 gold standard (PSG) and the scoring of the MB and the SC.

	PSG gold standard		
	WAKE	LIGHT SLEEP	DEEP SLEEP
MiBand (MB) staging			
Wake			
% Sensitivity	5.5	0.1	1.5
% PPV	62.8	4.7	32.6
Light sleep			
% Sensitivity	79.2	70.6	51.3
% PPV	18.9	57.8	23.2
Deep Sleep			
% Sensitivity	15.3	29.3	47.2
% PPV	7.5	48.9	43.6
Sleep Cycle (SC) staging			
Wake			
% Sensitivity	55.6	37.0	16.9
% PPV	24.3	61.1	14.7
Light sleep			
% Sensitivity	36.4	40.9	31.1
% PPV	14.4	61.2	24.4

Deep Sleep			
% Sensitivity	8.0	22.1	52.0
% PPV	4.1	42.8	53.0
Devices/applications	OA (%)		K
mi-band			
MB	53.31		0.14
Sleep Cycle			
SC	46.34		0.18

230 The agreement is demonstrated by the means of the sensitivity (%) as well as the positive predictive value
 231 (PPV). The percentage of the overall agreement (% OA) as well as the Cohen's Kappa coefficient (K) is reported
 232 for each device.

233 **Table 2.** Percentages of agreement table for 3 stages sleep scoring (Awake/ Asleep) between the gold standard
 234 (PSG) and the scoring of the MB and the SC.

	PSG gold standard	
	WAKE	SLEEP
Mi Band (MB) staging		
Wake		
% Sensitivity	5.5	0.5
% PPV	62.8	37.2
Sleep		
% Sensitivity	94.5	99.5
% PPV	13.2	86.8

Sleep Cycle (SC) staging

Wake		
% Sensitivity	55.6	32.6
% PPV	19.9	80.1
Sleep		
% Sensitivity	44.4	67.4
% PPV	8.7	91.3

MotionWatch (MW) staging

Wake		
% Sensitivity	37.5	7.8
% PPV	47.8	52.2
Sleep		
% Sensitivity	62.5	92.9
% PPV	11.5	88.5

Device/Application	OA (%)	K
Mi Band		
MB	86.54	0.08
Sleep Cycle		
SC	65.90	0.13
MotionWatch		
MW	83.42	0.33

235 4. Discussion

236 In the present study we evaluated the ability of 2 readily used consumer devices and one
 237 application for detecting true wake and sleep epochs at night. We acquired data from a commercial
 238 activity tracker, the Xiaomi Mi Band v2 (Xiaomi, Beijing, China), 2) a scientific actigraph, the
 239 CamNTEch Motionwatch 8 (CamNTEch, Cambridge, UK) and a readily used mobile phone
 240 application for sleep assessment, the Sleep Cycle App (v3.0.1.2511-release; Northcube, Gutenberg,
 241 Sweden). We then compared these consumer devices to our PSG gold standard which was
 242 simultaneously recorded. Overall, we revealed that these devices have an alarmingly low accuracy

243 in scoring sleep in three categories (that is, wake, light and deep sleep) with overall agreements
244 between 46.34% for the SC application and 53.02% for the wrist worn MB. If we tested only for the
245 correct classification in two classes, that is wake and sleep, the devices of course performed better
246 with an overall agreement of 65.90% for the SC, 84.69% for the MB and 81.33% for the scientific MW
247 device. Kappa coefficients however indicate only poor agreement with the PSG gold standard for
248 the MB and the SC (0.172 and 0.186 respectively) indicating serious bias in the scoring algorithms of
249 the MB device and the SC app.

250 We showed that all devices and applications had high accuracy in estimating the most global
251 sleep parameter, namely TiB. This makes these devices a helpful tool for objectively measuring the
252 time spent in bed at home rather than relying solely on subjective measures such as daily sleep
253 diaries. Especially in the case of the MW we need to note that we adjusted the start and the end of
254 the recordings to the PSG gold standard (as usually done using additional sleep diaries) which
255 might overestimate the fidelity of the MW device in measuring TiB. Nevertheless, our MW results
256 are consistent with those reported in previous literature [8,9].

257 Only for TiB correlational analysis also showed significant positive correlations between the
258 gold standard and all 3 sleep trackers. This raises the question of whether the faulty estimation of
259 values such as TST, SE, WASO or SOL are due to a priori knowledge of these sleep trackers of the
260 amount of time the average person actually sleeps or needs to fall asleep. If such information is
261 included in the algorithms and outputs of the consumer devices, this would explain why the largest
262 errors occur primarily for “non-average” sleep profiles and nights. However, to date this argument
263 remains speculative as all tested devices do not allow raw data access or are black boxes when it
264 comes to their staging algorithms. Similarly, when comparing the agreement between the MB and
265 SC with the PSG gold standard for 3 sleep-wake classes (light sleep, deep sleep, and wake) as
266 compared to 2 classes (sleep vs. wake) we found the expected increase in the OA yet a drop in the
267 Kappa scores. Especially for the MB, looking at the sensitivity scores we observed extremely low
268 sensitivity in detecting wakefulness (5.5%) and a very high sensitivity in detecting sleep (99.5%).
269 That is, by assigning “sleep” to basically every epoch the device also cannot miss sleep epochs, yet
270 it of course strongly overestimates sleep and has a vast amount of false alarms for stage “sleep”.
271 Although the MB was the least sensitive between all the 3 devices and applications, it had the
272 highest precision in scoring wakefulness (PPV: 62.8% for the MB, 47.8% for the MW and 24.3% for
273 the SC). That is, the MB does not score awakenings from sleep unless they are almost unmistakable.
274 This indicates a strong bias of the MB algorithm (as observed in the very low kappa values: 0.08;
275 [10]), which again raises the question if such a biased output can be of any benefit to the end-user.

276 Regarding the other key parameters evaluated, our results raise serious doubts whether such
277 consumer devices and applications can to-date provide any reliable information about sleep-related
278 health issues. Especially the revealed misjudgment in estimating key features of sleep such as SE,
279 SOL and WASO are worrisome as they are important diagnostic criteria for quantifying clinically
280 relevant bad sleep and sleep disorders such as insomnia [11]. On the contrary, by providing such
281 inaccurate information these consumer devices might even run risk to contribute to worse sleep and
282 life quality as end-users may be concerned by the sometimes negative output highlighting bad
283 nights of sleep [4].

284 Comparing the devices and SC app, our results suggest that wrist worn devices (MB and MW)
285 tend to have better a performance than mobile phone applications (SC) in measuring the key

286 features of sleep. This might be attributed to the fact that these devices have direct contact with the
287 body and hence are more accurate in capturing changes in physiological activity that accompanies
288 sleep of an individual and are therefore also more resilient to the environmental factors such as
289 noise or movement from the bed partner, child, pet or loud neighbors.

290 One very important drawback of the sleep trackers not mentioned yet is their inability to
291 provide any information about REM or “dreaming” sleep. Due to an inherent absence of needed
292 measurements for quantifying REM sleep (that is, most importantly eye movements via EOG and
293 brain activity via EEG) the devices on the market today cannot provide the full spectrum of sleep
294 even if the algorithms and sensors would be considerably improved. The incorporation of
295 additional sensors such as an eye or brain electrode might add substantially to the ability of these
296 devices to track and score sleep more accurately and in the long-run similar to a professional
297 polysomnography in the sleep laboratory.

298 An inherent limitation of our evaluation study is that most of our analysis need to build upon
299 the simple (graphical) outputs of the devices in form of plots provided for the end-user. For the
300 tested MB device and SC app there is no way to access the raw data. We therefore needed to come
301 up with a way to quantify the data and extract information that can be analyzed statistically (for
302 details see suppl. material). The MW device is a scientific device out of the price range of the usual
303 consumer and allows raw data access. Interestingly this device is likely the most accurate device
304 tested and yet its software only provides two outputs, sleep and wake classes, as it does not claim
305 to be able to classify sleep more fine-grained than that (as compared to consumer devices including
306 the MB and SC).

307 In summary, the currently available consumer devices for sleep tracking do not provide reliable
308 information about one’s sleep. However, devices of that kind could be very promising tools for
309 tracking sleep outside the laboratory in the future given that they adhere more to the scientific
310 standards of sleep staging and analysis. Moreover, by refining their algorithms or even adding
311 sensors these devices might be able to reliably monitor and classify sleep across its full range from
312 wakefulness to light sleep, deep sleep, and “REM” dreaming sleep.

313 **Author Contributions:** conceptualization, M.S. and M.A.; methodology, M.A. and M.S.; software, L.C. and
314 M.A.; validation, M.A. and M.S.; formal analysis, M.A. and L.C.; investigation, M.A, L.C., T.H. and M.H.
315 resources, M.S.; data curation, M.A. and L.C.; writing—original draft preparation, M.A.; writing—review and
316 editing, M.S. and M.A; visualization, L.C., M.A. and M.S.; supervision, M.S.; project administration, M.S.;
317 funding acquisition, M.S.”

318 **Supplementary Materials:** The following are available online at www.mdpi.com/xxx/s1. Figure S1: Sleep
319 scoring procedure of the sleep cycle application output hypnogram. Figure S2: Scatter plot showing the
320 positive correlations between the time in bed measurement of the PSG gold standard and those of the Mi Band,
321 MotionWatch and Sleep Cycle. Table S1: Spearman correlation results between the PSG gold standard and the
322 Mi Band, the MotionWatch and the Sleep Cycle application for the key sleep parameters. Table S2: Epoch-wise
323 agreement (Sleep/wake) while discarding epochs scored as REM by the PSG gold standard. Table S3:
324 Epoch-wise agreement comparison between 3 stage scoring (Wake/light Sleep/Deep sleep) and 2-stage scoring
325 (Sleep/wake).

326 **Funding:** This research received no external funding.

327 **Acknowledgments:** the authors would like to thank Ass.-Prof. Dr. rer. nat. Kerstin Hödlmoser, Monika
328 Angerer and Frank van Schalkwijk for their support throughout the whole process.

329 **Conflicts of Interest:** The authors declare no conflict of interest.

330 References

- 331 1. Portier, F., Portmann, A., Czernichow, P., Vascaut, L., Devin, E., Benhamou, D., ... Muir, J. F. (2000).
332 Evaluation of Home versus Laboratory Polysomnography in the Diagnosis of Sleep Apnea
333 Syndrome. *American Journal of Respiratory and Critical Care Medicine*, 162(3), 814–818.
334 <https://doi.org/10.1164/ajrccm.162.3.9908002>
- 335 2. Bruyneel, M., Sanida, C., Art, G., Libert, W., Cuvelier, L., Paesmans, M., ... Ninane, V. (2011). Sleep
336 efficiency during sleep studies: results of a prospective study comparing home-based and
337 in-hospital polysomnography. *Journal of Sleep Research*, 20(1pt2), 201–206.
338 <https://doi.org/10.1111/j.1365-2869.2010.00859.x>
- 339 3. Edinger, J. D., Fins, A. I., Sullivan, R. J., Marsh, G. R., Dailey, D. S., Hope, T. V., ... Vasilas, D. (1997).
340 Sleep in the Laboratory and Sleep at Home: Comparisons of Older Insomniacs and Normal Sleepers.
341 *Sleep*, 20(12), 1119–1126. <https://doi.org/10.1093/sleep/20.12.1119>
- 342 4. Baron, K. G., Abbott, S., Jao, N., Manalo, N., & Mullen, R. (2017). Orthosomnia: Are Some Patients
343 Taking the Quantified Self Too Far? *Journal of Clinical Sleep Medicine: JCSM: Official Publication of*
344 *the American Academy of Sleep Medicine*, 13(2), 351–354. <https://doi.org/10.5664/jcsm.6472>
- 345 5. Anderer, P., Gruber, G., Parapatics, S., Woertz, M., Miazhyńska, T., Klösch, G., ... Dorffner, G.
346 (2005). An E-Health Solution for Automatic Sleep Classification according to Rechtschaffen and
347 Kales: Validation Study of the Somnolyzer 24 × 7 Utilizing the Siesta Database.
348 *Neuropsychobiology*, 51(3), 115–133. <https://doi.org/10.1159/000085205>
- 349 6. Anderer, P., Moreau, A., Woertz, M., Ross, M., Gruber, G., Parapatics, S., ... Dorffner, G. (2010).
350 Computer-assisted sleep classification according to the standard of the American Academy of Sleep
351 Medicine: validation study of the AASM version of the Somnolyzer 24 × 7. *Neuropsychobiology*,
352 62(4), 250–264. <https://doi.org/10.1159/000320864>
- 353 7. American Academy of Sleep Medicine. (2007). The AASM manual for the scoring of sleep and
354 associated events: rules, terminology and technical specifications. Westchester, IL: American
355 Academy of Sleep Medicine, 23.
- 356 8. Kushida, C. A., Chang, A., Gadkary, C., Guilleminault, C., Carrillo, O., & Dement, W. C. (2001).
357 Comparison of actigraphic, polysomnographic, and subjective assessment of sleep parameters in
358 sleep-disordered patients. *Sleep Medicine*, 2(5), 389–396.
359 [https://doi.org/10.1016/S1389-9457\(00\)00098-8](https://doi.org/10.1016/S1389-9457(00)00098-8)
- 360 9. Sadeh, A. (2011). The role and validity of actigraphy in sleep medicine: An update. *Sleep Medicine*
361 *Reviews*, 15(4), 259–267. <https://doi.org/10.1016/j.smrv.2010.10.001>
- 362 10. Flight, L., & Julious, S. A. (2015). The disagreeable behaviour of the kappa statistic. *Pharmaceutical*
363 *Statistics*, 14(1), 74–78. <https://doi.org/10.1002/pst.1659>
- 364 11. Schutte-Rodin, S., Broch, L., Buysse, D., Dorsey, C., & Sateia, M. (2008). Clinical Guideline for the
365 Evaluation and Management of Chronic Insomnia in Adults. *Journal of Clinical Sleep Medicine :*
366 *JCSM: Official Publication of the American Academy of Sleep Medicine*, 4(5), 487–504.
367