
DreamCatcher: A Wearer-aware Sleep Event Dataset Based on Earables in Non-restrictive Environments

Zeyu Wang^{1*}, Xiyuxing Zhang^{1*}, Ruotong Yu^{1*};
Yuntao Wang^{1†}, Kenneth Christofferson², Jingru Zhang¹, Alex Mariakakis²,
Yuanchun Shi^{1,3}

¹ Department of Computer Science and Technology, Tsinghua University[‡]

² University of Toronto ³ Qinghai University
wang-zy23@mails.tsinghua.edu.cn, yuntaowang@tsinghua.edu.cn

Abstract

1 Poor quality sleep can be characterized by the occurrence of events ranging from
2 body movement to breathing impairment. Widely available earbuds equipped
3 with sensors (also known as earables) can be combined with a sleep event de-
4 tection algorithm to offer a convenient alternative to laborious clinical tests for
5 individuals suffering from sleep disorders. Although various solutions utilizing
6 such devices have been proposed to detect sleep events, they ignore the fact that
7 individuals often share sleeping spaces with roommates or couples. To address
8 this issue, we introduce DreamCatcher, the first publicly available dataset for
9 wearer-aware sleep event algorithm development on earables. DreamCatcher
10 encompasses eight distinct sleep events, including synchronous dual-channel au-
11 dio and motion data collected from 12 pairs (24 participants) totaling 210 hours
12 (420 hour.person) with fine-grained label. We tested multiple benchmark mod-
13 els on three tasks related to sleep event detection, demonstrating the usability
14 and unique challenge of DreamCatcher. We hope that the proposed Dream-
15 Catcher can inspire other researchers to further explore efficient wearer-aware
16 human vocal activity sensing on earables. DreamCatcher is publicly available at
17 <https://github.com/thuhci/DreamCatcher>.

18 1 Introduction

19 More than one-seventh of the global population suffers from at least one kind of sleep disorder, yet
20 many are undiagnosed [6, 36, 41]. Sleep disorders can lead to various health issues, such as cardiovas-
21 cular disease and depression [14, 20, 39]. The gold-standard diagnostic method, polysomnography
22 (PSG), requires patients to spend the night in a specialized sleep clinic. Conducting such sleep
23 studies can be cost-prohibitive and resource-intensive. Additionally, patients may suffer from the
24 “first-night effect” where they exhibit anomalous sleep behavior when spending the night in a new

*equal contribution

†corresponding author

‡Key Laboratory of Pervasive Computing, Ministry of Education, Beijing National Research Center for Information Science and Technology, Department of Computer Science and Technology, Tsinghua University, Beijing, 100084, China

25 environment [29]. These challenges call for a minimally intrusive at-home sleep monitoring solution
26 that can alert wearers to potential sleep disorders.

27 Many sleep disorders are associated with at least one detectable sleep event. For instance, obstructive
28 sleep apnea (OSA) is characterized by the sudden cessation of snoring [4, 30], bruxism manifests
29 as frequent teeth grinding or clenching [19], and restless sleep is often accompanied by excessive
30 nighttime movement [41]. Because it is difficult for people to recall these events while sleeping,
31 continuous monitoring is crucial to facilitate diagnosis.

32 Recent research has shown that lightweight earables [8] can provide convenient real-time monitoring
33 of human activity [11, 34, 45, 52, 46]. For sleep monitoring in particular, earables have unique
34 advantages over other wearables like smartwatches and smartphones [3, 9, 24, 25, 48]. The ears
35 are located on the head and close to the trunk of the body, allowing microphones to capture rich
36 acoustic information generated during sleep. The in-ear feedback microphone included in active
37 noise-cancelling earbuds can even detect subtle sounds produced within the body. For this work, we
38 utilize a modified commercial earbud containing two microphones (feedback and feedforward) and an
39 inertial measurement unit (IMU).

40 Advancements in hardware technology and machine learning algorithms have spurred increased
41 research into sleep monitoring using commodity wearables. Current acoustic-based sleep event
42 detection algorithms mainly focus on audio feature engineering [2, 13] or lightweight deep learning
43 models [11]. These solutions are often developed using data collected in controlled environments and
44 contrived scenarios with minimal confounds (e.g., ambient noise). However, people often share sleep
45 spaces with other individuals like roommates or spouses who may move and create sounds, leading
46 to observable events not associated with the wearer [9, 13]. Moreover, these studies have not made
47 their code or datasets publicly available.

48 We address these shortcomings by presenting and releasing DreamCatcher — a large-scale, multi-
49 modal, multi-sleeper sleep event dataset of earable data with fine-grained labels. We recruited 12 pairs
50 (24 participants) of people who slept in the same room, and one person from each pair had a potential
51 sleep disorder. We collected earable data from these pairs over the course of 420 hours and manually
52 annotated 8 sleep events: teeth grinding, swallowing, somniloquy, breathing, coughing, snoring, and
53 body movement. To demonstrate the utility of DreamCatcher, we present case studies of how it
54 can be used to train baseline models that address three valuable tasks: wearer event identification,
55 wearer-aware sleep sound event classification, and wearer-aware sleep sound event detection.

56 The main contributions of this work are as follows:

- 57 • We collected and released the first and largest sleep dataset based on multi-modal earable data
58 collected in real scenarios with the disruption of sleep partners. Data is synchronized and
59 annotated with fine-grained event labels.
- 60 • We benchmarked DreamCatcher on three sleep monitoring tasks: wearer event identification,
61 wearer-aware sleep sound event classification, and wearer-aware sleep sound event detection.
- 62 • We provide open-source resources including the dataset, code for setting up benchmarks, and
63 tutorial for constructing the earable hardware we used.

64 **2 Related Work**

65 **2.1 Contactless Sleep Monitoring and Wearer-Awareness**

66 The gold standard for sleep monitoring is polysomnography (PSG), which entails wiring a series of
67 sensors onto an individual in a sleep clinic for continuous monitoring and observation. PSG sessions
68 are expensive, labour-intensive, and time-consuming [28], so researchers have shown substantial
69 interest in developing more convenient sleep monitoring solutions suitable for home use.

70 Contactless sleep monitoring typically falls under one of two methods. The first method involves
71 acoustic sensing of audible sounds using smartphones [24], smartwatches [9, 12], and earbuds [2, 11,

72 13, 33]. While these systems can work with commodity devices, they are prone to interference in
 73 multi-user settings. The second method relies on wireless sensing to detect body motion, respiration,
 74 and even heartbeats through minor chest movements at specific frequencies. Commonly used signals
 75 include WiFi [50], mmWave [49], and sonar [27]. Wireless sensing makes it possible to manage
 76 multi-user scenarios since reflections from multiple users arrive at different times [49, 27]. However,
 77 dedicated devices for such approaches are non-trivial to deploy, and wireless sensing is less effective
 78 for detecting sleep events such as snoring, swallowing, or somniloquy.

79 In summary, acoustic and wireless solutions to contactless sleep monitoring show promise in address-
 80 ing the multi-user challenge, yet cater to different aspects of sleep monitoring. Moreover, as indicated
 81 in [15], there is an inherent trade-off between accuracy and comfort.

82 2.2 Sleep Monitoring with Earables

83 Compared to other commodity wearables for sleep monitoring, earables are worn in closer proximity
 84 to respiratory-vocal system and the external carotid artery, offering an ideal position for measuring
 85 behaviors and physiological parameters related to sleep [35]. These opportunities have been leveraged
 86 using specialized biomedical sensors for sleep monitoring around the ear, such as in-ear EEG [23, 26,
 87 40] and PPG [44] sensors, but these sensors are not widely available on commercial earables due to
 88 their high cost and integration complexity.

89 Some recent works have explored sleep monitoring by leveraging earables without modification,
 90 relying on motion sensors and in-ear microphones used for active noise cancellation. Leveraging
 91 the audio signals from earables, Ren et al. [33] developed a system that could track breathing rate
 92 and detect four sleep events. Christofferson et al. [11] utilized microphones in commercial earbuds
 93 for sleep sound classification. Their proposed SleepTSM model achieved promising performance
 94 in detecting seven different sleep events with a small footprint suitable for deployment on earables.
 95 Han et al. [13] proposed EarSleep, a similar sleep stage classification system dependent on acoustic
 96 sensing of body sounds.

97 Although the microphones on commercial earables have been used to great effect in sleep monitoring,
 98 such systems are often evaluated in controlled scenarios with a single participant at a time. In
 99 multi-sleeper scenarios, sounds may originate from people who are not wearing the earable, leading
 100 to mischaracterizations of the wearer’s sleep experiences. Drawing inspiration from the EarSAVAS
 101 dataset [51], our work on DreamCatcher facilitates the development and evaluation of wearer-
 102 aware ubiquitous acoustic sleep event monitoring systems by providing a public sleep dataset that
 103 encompasses not only the wearer’s sound events but also interference from non-wearers and non-
 104 restrictive environmental conditions.

105 2.3 Sleep Datasets

106 Table 1 compares datasets across sensing modalities that have been leveraged for sleep monitoring
 107 research. It reveals multiple data-related challenges faced by previous works:

Method	Device	Modalities		Scale		Scenario		Open-source
		Acoustic	Non-Acoustic	Data Amount	Participants	Real	Open	
SleepEDF [17]	PSG	✗	✓	197 nights	–	✓	✗	✓
MASS [31]	PSG	✗	✓	200 nights	200	✓	✗	✓
SleepHunter [12]	smartphone	✓	✓	90 nights	45	✓	✗	✗
SleepGuard [9]	smartwatch	✓	✓	210 nights	15	✓	✗	✗
FusedTSMNet [2]	–	✗	✗	1 hour	–	✗	✗	✗
SleepTSM [11]	earable	✓	✗	6 hours	20	✗	✗	✗
EarSleep [13]	earable	✓	✗	48 nights	18	✓	✗	✗
Ren et al. [33]	earable/smartphone	✓	✗	–	6	✓	✗	✗
DreamCatcher (ours)	earable	✓ (2-channel)	✓	420 hours / 62 nights	24	✓	✓	✓

Real: whether events were real or simulated; Open: whether multiple individuals were in the same room

Table 1: Sleep Study Dataset Comparison.

- 108 1. Although there has been substantial work utilizing ubiquitous sensors such as microphones in
109 commodity wearables, the only publicly available datasets are for gold-standard PSG.
- 110 2. Previous wearable solutions reliant on audio phenomena have overlooked the interference of
111 non-wearer in multiple sleeper scenarios.
- 112 3. Proprietary datasets using commodity wearables are typically limited in scale, both in terms
113 of the number of participants and the quantity of data per person. This issue is particularly
114 prevalent in research related to earables [34].

115 Our dataset aims to fill the gaps. To the best of our knowledge, DreamCatcher is the first open-source
116 sleep event dataset targeted at ubiquitous sensors on commercial devices. Previous work has only
117 considered single-sleeper scenarios. By integrating data from non-wearers, DreamCatcher facilitates
118 the development and evaluation of wearer-aware sleep event monitoring. Our dataset consists of
119 synchronous dual-channel audio and motion data collected from 12 pairs (24 participants) totaling
120 210 hours (420 hour.person) with fine-grained labels of eight distinct sleep events. As the largest
121 open-source sleep dataset to date, we envision DreamCatcher will advance wearer-aware sleep event
122 monitoring on commercial earables.

123 3 DreamCatcher Dataset

124 3.1 Dataset Collection

125 **Hardware.** Using a commodity earable is important because custom devices can often be optimized
126 for data quality in ways that do not translate to existing platforms. Because commodity earables
127 do not provide API access to their data streams, we had to modify an earbud for data acquisition.
128 As shown in Figure 1a, we integrated an MPU6050 IMU sensor into the hardware of Bose QC 20
129 earbuds, preserving the native feedback and feedforward microphone configuration. All sensors
130 were controlled by a compact external development board, wherein the audio signal was sampled
131 at 24 kHz and the IMU signal was sampled at approximately 94 Hz. To enhance user comfort, the
132 development board was integrated into an enclosure and wrapped like a necklace. This device can
133 function continuously for roughly 7 hours. Appendix A.1 contains more implementation details.

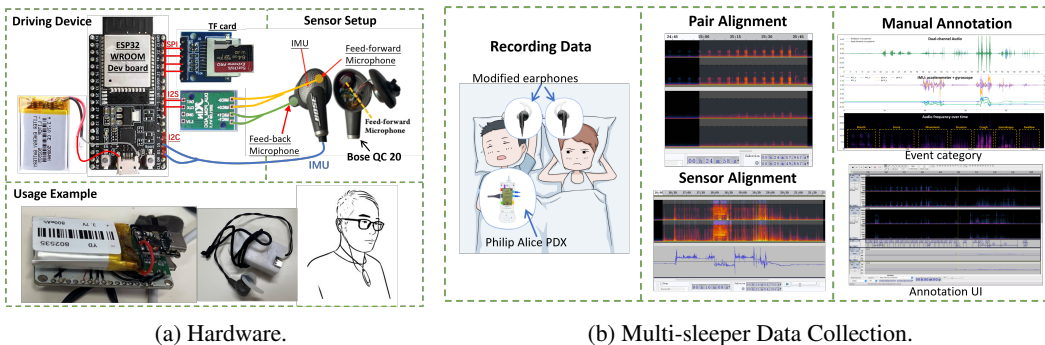


Figure 1: Experiment Setup.

134 For gold-standard sleep monitoring data, we used a portable PSG system by Philips called the
135 Alice PDX⁴. This device includes a canula and thermal sensor for measuring airflow in the nose, a
136 chest strap for measuring chest expansion during breathing, a fingertip SpO2 sensor for measuring
137 peripheral oxygen saturation, and limb movement sensors. This device was only worn by the person
138 in each participant pair who reported a sleep disorder.

⁴<https://www.usa.philips.com/healthcare/product/HC1043844/alice-pdx-portable-sleep-diagnostic-system>

139 **Participants.** To collect data while accounting for the influence of a sleep partner, we recruited
140 participants in dyads, imposing no constraints on their relationship as long as they shared a bedroom.
141 A total of 12 pairs (24 participants) participated in the study, including 9 males and 15 females aged
142 between 19 and 51 (average = 24.7, standard deviation = 8.3). Among the 12 participants who were
143 equipped with the PDx device, 6 individuals were observed to exhibit varying degrees of sleep apnea.

144 **Experiment Protocol.** As shown in the first step of Figure 1b, each pair of participants slept in the
145 same bedroom for at least 6 hours. They were instructed to start their sleep session around the same
146 time to maximize the amount of temporal overlap in their data. They were also asked to set an alarm
147 clock that could be heard by both earphones; this sound was used for post-hoc manual data alignment.
148 After the participants woke up and turned off the data collection hardware, they were required to fill
149 out a PSQI [7] questionnaire to self-report their sleep quality.

150 3.2 Annotation and Statistics

151 **Data Alignment.** As shown in the second step of Figure 1b, we performed post-hoc alignment
152 for (1) the audio and IMU data in each earable and (2) each pair’s audio data. The first round of
153 alignment involving the data modalities within each earbud would hypothetically be trivial because
154 the sensors should be intrinsically synchronized as they are connected to the same ESP32 board and
155 controlled by the same microcontroller. However, we observed an inherent clock drift between the
156 audio and IMU sampling protocols. Over a span of 7 hours, the IMU recording extended 3 seconds
157 longer than the audio, accounting for a deviation of about 0.01%. To correct this, we re-scale the
158 IMU data to match the audio recording duration, as the drift is evenly distributed over the entire
159 recording period.

160 Because each participant’s data was recorded independently by separate earbuds, the second round of
161 alignment involved aligning data across participant pairs. To accomplish this, we utilized an alarm
162 clock as a compensatory reference, manually adjusting the audio recordings to align with the alarm
163 clock’s spectrogram.

164 **Annotation.** Because data was collected from participants’ homes, using video was not an accept-
165 able form of annotation due to privacy concerns. Instead, we set up a hierarchical inspection process
166 in which a team of annotators reviewed the earbud data to identify and label events. The annotators
167 were asked to identify the eight sleep events listed in Table 2. They used Audacity⁵ to inspect each
168 participant’s binaural audio channel and IMU data as well as the sleep partner’s binaural audio data
169 simultaneously. The IMU data helped annotators determine the category of wearer-emitted events,
170 while the sleep partner’s audio helped them determine whether the event was emitted by the wearer.
171 The annotation process, described more thoroughly in Appendix A.4, entailed selecting an interval for
172 each event they noticed and then assigning a category to it. Each label was checked by at least three
173 annotators; whenever they did not reach a consensus, voting was used to assign labels. Examples of
174 each event are provided in Figure 2.

175 **Dataset Statistics.** Table 2 summarizes the prevalence and duration of each event type, while
176 Figure 3a illustrates the distribution of durations of each event type. DreamCatcher is a highly
177 imbalanced dataset, reflecting the natural scarcity of certain sleep disturbances such as bruxism,
178 swallowing, somniloquy, and coughing.

179 Figure 3b shows the smoothed average frequency of different events over the course of a typical
180 night of sleep. Note that participants slept for different amounts of time, so there may be some
181 misalignment in the timing of events across individuals. However, the plot reflects some known
182 observations about sleeping. For example, movement and swallowing were less prevalent after the
183 first hour of sleep, while snoring and somniloquy became more prevalent.

⁵<https://www.audacityteam.org/>

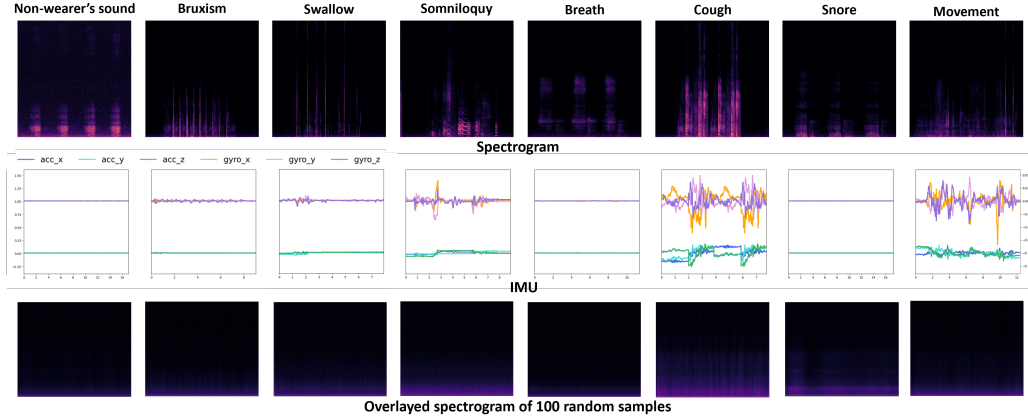


Figure 2: Examples of Each Sleep Event.

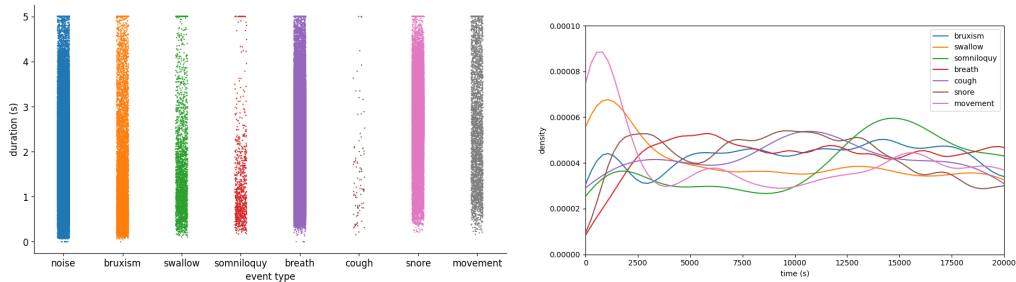
Label	Description	Total Duration (hrs)	Avg Duration (secs)	S.D. (secs)
Noise	Acoustic events emitted by non-wearers, as well as background noises	32.78	2.27	2.49
Bruxism	Grinding or clenching teeth	5.15	2.00	4.01
Swallow	Reflexively or intentionally saliva swallowing	1.28	1.75	1.89
Somniloquy	Talking aloud, murmuring, or shouting while asleep	0.37	1.49	2.39
Breathe	One inhalation + one exhalation	83.56	2.21	2.41
Cough	Coughing, throat clearing, or sniffing	0.04	1.60	1.24
Snore	One inhalation + one exhalation with prominent vibrations or whistling	31.98	3.10	3.97
Movement	Shifts in position or gestures made	10.51	6.83	7.17

Table 2: Label Definitions and Summary Statistics.

184 **Participant Data Splits.** To evaluate how a sleep monitoring system would generalize to unseen
 185 users without any calibration data, we recommend splitting data according to participant IDs. Each
 186 participant in our dataset exhibited different distributions of sleep events; this information is detailed
 187 in Appendix A.4 A.2. The split configurations that optimize the balance in label prevalence are
 188 depicted in Figure 4.

189 3.3 Ethics and Accessibility

190 The protocol used to generate the DreamCatcher dataset received approval from the local Institutional
 191 Review Board (IRB) where the data was collected. Participants were explicitly informed about
 192 the data recording process and that the dataset would be made publicly available. To safeguard
 193 participants' privacy, DreamCatcher has been fully anonymized. An important consideration in
 194 this regard is the fact that participants recorded data at home for an entire night, so any private



(a) Label duration distribution for each category.

(b) Averaged label density overnight.

Figure 3: Label Distributions.

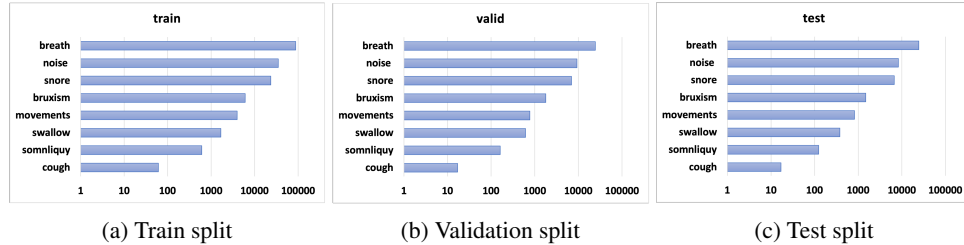


Figure 4: Cross-user Splits.

195 conversations may have been recorded by the earables’ microphones. To address this concern, all
 196 non-somniquy dialogue was manually removed from the dataset before it was released.

197 4 Benchmarks

198 4.1 Wearer Event Identification

199 **Task Description.** In multi-sleeper scenarios, sleep event monitoring using acoustic methods
 200 suffers from the interference of external sounds not produced by the earable wearer. Therefore, we
 201 define wearer identification as a binary classification task focused on determining whether audio
 202 events come from the wearer or other sources. To the best of our knowledge, no previous work has
 203 explored this topic in the context of sleep event monitoring.

204 **Dataset Preparation.** For this task, we assume that candidate events have been separated from
 205 silence using a simple threshold-based approach and focus only on data that our annotators labeled.
 206 After segmenting synchronous dual-channel audio and motion data according to the fine-grained
 207 labels, we extracted the features listed in Table 3a from each event. The primary challenge for this
 208 task is the design of features that are computationally efficient for wearable devices with limited
 209 processing capabilities, capable of distinguishing events caused by the sleep partner.

210 Due to the low intensity of background noise typically seen in real-world sleep scenarios, we relied
 211 on acoustic features to distinguish between wearer and non-wearer events. We calculated traditional
 212 acoustic features like zero-crossing rate (ZCR) and root mean square (RMS) on both the feedforward
 213 and feedback channels. We also calculated three inter-channel audio features — RMS-ED, Mel-FD,
 214 and TDOA — that model the different propagation characteristics of sounds. Based on the observation
 215 that bone-conducted sound from the wearer should have higher energy at the feedback microphone
 216 than at the feedforward microphone, RMS-ED measures the root mean square of the energy difference
 217 between the two audio channels. Given the different propagation paths of wearer and non-wearer
 218 sounds reaching the ear, Mel-FD measures the energy difference between the Fast Fourier Transforms
 219 (FFTs) calculated from both audio channels according to the Mel scale. Finally, time difference of
 220 arrival (TDOA) between the two channels reflects the propagation path difference between wearer
 221 and non-wearer sounds.

222 Since the wearer’s sleep events are often accompanied by body movements that are captured by the
 223 earables’ motion sensors, we also extracted motion-related features. We first calculated the overall
 224 magnitude of the accelerometer and gyroscope data separately, after which we computed IMU-STD
 225 as the mean and standard deviation of those magnitudes over time.

226 **Benchmark Methods.** Given the low dimensionality of the input data, we benchmarked five
 227 traditional machine learning models: two low-complexity models (logistic regression and linear
 228 SVM) and three high-complexity models (random forest, decision tree, and AdaBoost).

229 **Model Training and Evaluation Metrics.** To evaluate the performance of the models, we trained
 230 and evaluated each one using leave-one-user-out cross-validation. Since this is a binary classification
 231 task, we report performance according to accuracy, F1 score, precision, recall, and AUC.

232 **Results.** Table 3a shows that most of our features could be computed using fewer than 1 M FLOPs
 233 across our entire dataset. The results in Table 3b demonstrate that the high-complexity models
 234 achieved similarly higher accuracy compared to the low-complexity ones. According to the feature
 235 importance scores presented in Appendix B.1, we found that models with higher complexity are more
 236 effective in leveraging the inter-channel audio features. Furthermore, we observed that motion features
 237 were important for all models, particularly those that were less complex. These results highlight the
 238 utility of inter-channel audio and motion features for wearer awareness of sleep monitoring.

Modality		FLOPs (M)	Method	Acc.	AUC	Recall	Prec.	F1
ZCR	per-channel (audio)	0.4	Random Forest	0.997	0.999	0.998	0.998	0.998
RMS	per-channel (audio)	0.4	Decision Tree	0.993	0.990	0.996	0.996	0.996
RMS-ED	inter-channel (audio)	0.6	AdaBoost	0.922	0.965	0.960	0.947	0.951
Mel-FD	inter-channel (audio)	6.1	Logistic Regression	0.797	0.597	0.994	0.800	0.886
TDOA	inter-channel (audio)	2.7	SVM (linear)	0.602	0.539	0.645	0.816	0.721
IMU-STD	per-sensor (accel & gyro)	0.5						

(a) Comparison of Extracted Features.

(b) Comparison of ML Algorithms.

Table 3: Benchmarks for Wearer Event Identification.

239 4.2 Wearer-Aware Sleep Sound Event Classification

240 **Task Description.** Sleep sound event classification serves as the foundation of sleep disorder
 241 diagnosis. Han et al. [13] also revealed that the categorization of sleep sound events also facilitates
 242 sleep stage inference. Although algorithms already exist for this task, the interference caused by
 243 non-wearers is often overlooked, limiting their applicability in multi-sleeper scenarios. Inspired
 244 by EarSAVAS [51], we define wearer-aware sleep sound event classification as an $(n + 1)$ -class
 245 multi-classification task, where n represents the number of target events and the remaining class
 246 encompasses both ambient and non-wearer sounds.

247 **Dataset Preparation.** As with wearer event identification, we assume that event onset and offset
 248 are already known for this task. To standardize the input data size, we cropped the synchronous audio
 249 and motion data into 5-second clips that were sufficiently long to cover the duration of the longest
 250 event in our dataset.

251 **Benchmark Methods.** We examined five state-of-the-art models for this task:

- 252 1. **SleepTSM** [11] is a lightweight sleep sound classification model that was not evaluated with
 253 multiple sleepers in the same room.
- 254 2. **EarVAS** [51] and its variants were evaluated on the EarSAVAS dataset to demonstrate subject-
 255 aware vocal activity classification utilizing dual-channel audio and motion data.
- 256 3. **Wav2Vec2.0** [5], **BEATs** [10], and **CLAP** [47] are generic audio event classification methods.

257 Besides EarVAS, the other models are only designed to support single-channel audio input. Since
 258 DreamCatcher includes audio from both the feedforward and feedback microphones, we evaluated
 259 these models on each of those channels separately. All model pre-processing steps and hyperparam-
 260 eters were configured identically to those in the original works we replicated. Appendix C.1 shows the
 261 details of the partition of our dataset and the training details of every benchmark model.

262 **Model Training and Evaluation Metrics.** Each model was evaluated using leave-one-user-out
 263 cross-validation. Since this is a multi-class task, we used accuracy, macro-averaged AUC, macro-
 264 averaged F1 score, and MCC as evaluation metrics. We also report model complexity according to
 265 FLOPs and the number of parameters.

Method	input channel	Evaluation Metrics (%)				FLOPs (G)	Params. (M)
		Acc.	Macro-AUC	Macro-F1	MCC		
SleepTSM [11]	feedback	73.01	63.51	35.61	49.98	0.927	0.37
SleepTSM	feedforward	72.89	64.00	36.97	50.76	0.927	0.37
EarVAS [51]	all	78.07	74.49	36.76	49.22	0.354	12.90
EarVAS	dual-channel-audio	76.64	77.36	38.77	51.12	0.040	4.40
EarVAS	feedback	75.75	75.75	34.68	46.84	0.040	4.40
EarVAS	feedforward	76.99	75.06	34.76	43.02	0.040	4.40
EarVAS	imu-only	45.03	60.56	13.75	12.24	0.313	8.50
BEATs [10]	feedback	90.73	80.38	57.47	66.60	22.46	90.51
BEATs	feedforward	89.64	78.93	55.51	59.04	22.46	90.51
Wav2Vec2.0 [5]	feedback	75.45	91.60	48.36	56.72	26.84	94.39
Wav2Vec2.0	feedforward	73.29	88.84	42.52	54.11	26.84	94.39
CLAP (zero-shot) [47]	feedback	37.04	65.57	16.77	18.21	53.03	190.80
CLAP (zero-shot)	feedforward	37.35	65.64	17.31	18.98	53.03	190.80

Table 4: Benchmarks for Wearer-Aware Sleep Sound Event Classification.

266 **Results.** As shown in Table 4, most of the models achieved accuracies above 70%; the exceptions
267 were CLAP and a configuration of EarVAS that only used IMU data. However, the macro-F1 scores
268 were typically far lower. This is largely due to the significant class imbalance of our dataset, as
269 some events are far more common than others. Another hurdle encountered by these models was the
270 challenge of jointly optimizing wearer event identification and sleep sound classification. We used a
271 single model to perform both tasks simultaneously, but a dual-stage pipeline may be more appropriate
272 in future work. Appendix C.2 provides a more thorough analysis of the results, showing the efficacy
273 of the feedback microphone in detecting low-intensity events like swallowing and highlighting the
274 promise of sensor fusion for future explorations.

275 4.3 Wearer-Aware Sleep Sound Event Detection

276 **Task Description.** Knowing when a sleep event starts and stops is crucial for sleep monitoring,
277 as the temporal distribution and order of events provide critical insights into sleep progression [13].
278 Inspired by sound event detection (SED) systems and the DCASE Challenge [18], we define wearer-
279 aware sleep sound event detection as a task that involves determining not only the category of an
280 event but also its onset and offset.

281 **Dataset Preparation.** Following the data format standards from the DCASE Challenges [18], we
282 used 10-second clips for this experiment so that the models would have enough context for precise
283 event detection.

284 **Benchmark Methods.** State-of-the-art sound event detection methods predominantly employ deep
285 learning, with most of them being built upon convolutional recurrent neural networks (CRNNs).
286 According to Mesaros et al. [22], such methods have been both trained from scratch and have utilized
287 transfer learning to shortcut learning. We benchmarked SEDNet [1] and ATST-SED [38] to represent
288 these two categories, respectively. We selected SEDNet because of its pioneering role in using
289 CRNNs with multi-channel microphone data for sound event detection. On the other hand, we
290 selected ATST-SED because it outperformed all competitors on the DESED dataset [42].

291 **Model Training and Evaluation Metrics.** Each model was evaluated using leave-one-user-out
292 cross-validation. We used conventional collar-based metrics [21] including event-based macro-
293 averaged F1 score and error rate to quantify model performance.

294 **Results.** According to the results shown in Table 5, we found that ATST-SED achieved significantly
295 better performance at the cost of a much larger footprint. We also observed that both models were
296 more accurate when they were trained using feedforward microphone audio. In fact, SEDNet trained
297 on multiple audio channels achieved the lowest macro-F1 score out of all the configurations we tested.
298 Appendix D.3 provides a more thorough analysis of the results.

Method	input channel	Evaluation Metrics		FLOPs (G)	Params. (M)
		Macro-F1 (%)	Error Rate		
SEDNet [1]	dual-channel audio	14.02	0.97	0.31	0.37
SEDNet	feedback	18.85	0.91	0.30	0.37
SEDNet	feedforward	17.98	0.98	0.30	0.37
ATST-SED [38]	feedback	24.73	0.85	44.16	172.9
ATST-SED	feedforward	24.10	0.85	44.16	172.9

Table 5: Benchmarks for Wearer-Aware Sleep Sound Event Detection.

5 Limitations and Future Work

First of all, the natural frequency distribution of sleep events leads to a highly imbalanced dataset in DreamCatcher, with rarer events often holding greater significance. Based on the DreamCatcher dataset, we generated a balanced dataset through data augmentation methods and trained classification models on it, as shown in Appendix C.3. We envision the generation of rare sleep events, which must aligns with the patterns of human sleep, will be a highly valuable area for future research.

Moreover, privacy concerns in multi-sleeper settings preclude video verification, resulting in potential label inaccuracies despite requiring a consensus among at least three annotators for challenging-to-identify events. Furthermore, the emergence of commercial earphones equipped with physiological sensors like photoplethysmography (PPG) presents an opportunity to enhance DreamCatcher with additional data modalities in future iterations.

The current prototype earbuds used in our study may not be the epitome of comfort for all users, especially for those who have difficulties falling asleep. However, the existence of commercially available sleep earbuds that are small, soft, and ergonomically designed (e.g., Amazfit Zenbuds⁶ and Bose Sleepbuds⁷) underscores the potential for earbuds to become a comfortable and viable sleep monitoring platform. These options point towards a promising future for the application of earbud technology in sleep studies.

6 Conclusion

This paper introduces DreamCatcher, the first open-source dataset featuring multi-sleeper, multi-modal data from a commodity device along with fine-grained annotations of sleep disorder-related sound events. DreamCatcher encompasses 420 hours of synchronized dual-channel audio and motion data, offering a rich and challenging resource for sleep monitoring. We validated DreamCatcher’s utility by establishing benchmarks across three distinct tasks, and we hope that these results motivate other researchers to innovate further on our dataset.

7 Acknowledgement

This work is supported by Natural Science Foundation of China under Grant No. 62472244, No. 62132010 and No. 62222606, University of Toronto – Tsinghua University Joint Research Fund, Natural Sciences and Engineering Research Council of Canada (NSERC) Discovery Grant (#RGPIN-2021-03457), Tsinghua University Initiative Scientific Research Program, Beijing Key Lab of Networked Multimedia, Institute for Artificial Intelligence, Tsinghua University (THUAI). Thank to the participants involved in data collection and to the professional annotators who contributed to data labeling.

⁶<https://www.amazfit.com/products/amazfit-zenbuds>

⁷<https://www.bose.ca/en/p/all-health/bose-noisemasking-sleepbuds/SBD-SLEEPBUDS.html>

References

- [1] S. Adavanne, P. Pertilä, and T. Virtanen. Sound event detection using spatial features and convolutional recurrent neural network. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 771–775, 2017.
- [2] E. Akbal and T. Tuncer. Fusedtsnet: An automated nocturnal sleep sound classification method based on a fused textural and statistical feature generation network. *Applied Acoustics*, 171:107559, 2021.
- [3] F. Al Hossain, A. A. Lover, G. A. Corey, N. G. Reich, and T. Rahman. Flusense: a contactless syndromic surveillance platform for influenza-like illness in hospital waiting areas. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 4(1):1–28, 2020.
- [4] A. M. Alencar, D. G. V. da Silva, C. B. Oliveira, A. P. Vieira, H. T. Moriya, and G. Lorenzi-Filho. Dynamics of snoring sounds and its connection with obstructive sleep apnea. *Physica A: Statistical Mechanics and its Applications*, 392(1):271–277, 2013.
- [5] A. Baevski, H. Zhou, A. Mohamed, and M. Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations, 2020.
- [6] A. V. Benjafield, N. T. Ayas, P. R. Eastwood, R. Heinzer, M. S. Ip, M. J. Morrell, C. M. Nunez, S. R. Patel, T. Penzel, J.-L. Pépin, et al. Estimation of the global prevalence and burden of obstructive sleep apnoea: a literature-based analysis. *The Lancet Respiratory Medicine*, 7(8):687–698, 2019.
- [7] D. J. Buysse, C. F. Reynolds III, T. H. Monk, S. R. Berman, and D. J. Kupfer. The pittsburgh sleep quality index: a new instrument for psychiatric practice and research. *Psychiatry research*, 28(2):193–213, 1989.
- [8] Canalsys. Global smart device shipment forecasts 2020 to 2023, 2020. Last accessed January 2020.
- [9] L. Chang, J. Lu, J. Wang, X. Chen, D. Fang, Z. Tang, P. Nurmi, and Z. Wang. Sleepguard: Capturing rich sleep information using smartwatch sensing data. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 2(3), sep 2018.
- [10] S. Chen, Y. Wu, C. Wang, S. Liu, D. Tompkins, Z. Chen, and F. Wei. Beats: Audio pre-training with acoustic tokenizers, 2022.
- [11] K. Christofferson, X. Chen, Z. Wang, A. Mariakakis, and Y. Wang. Sleep sound classification using anc-enabled earbuds. In *2022 IEEE International Conference on Pervasive Computing and Communications Workshops and other Affiliated Events (PerCom Workshops)*, pages 397–402, 2022.
- [12] W. Gu, L. Shanguan, Z. Yang, and Y. Liu. Sleep hunter: Towards fine grained sleep stage tracking with smartphones. *IEEE Transactions on Mobile Computing*, 15(6):1514–1527, 2015.
- [13] F. Han, P. Yang, Y. Feng, W. Jiang, Y. Zhang, and X.-Y. Li. Earsleep: In-ear acoustic-based physical and physiological activity recognition for sleep stage detection. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 8(2):1–31, 2024.
- [14] E. A. Hill. Obstructive sleep apnoea/hypopnoea syndrome in adults with down syndrome. *Breathe*, 12(4):e91–e96, 2016.
- [15] Z. Hussain, Q. Z. Sheng, W. E. Zhang, J. Ortiz, and S. Pouriyeh. Non-invasive techniques for monitoring different aspects of sleep: A comprehensive review. *ACM Transactions on Computing for Healthcare (HEALTH)*, 3(2):1–26, 2022.

- 374 [16] K. Imoto, S. Mishima, Y. Arai, and R. Kondo. Impact of data imbalance caused by inactive
375 frames and difference in sound duration on sound event detection performance. *Applied*
376 *Acoustics*, 196:108882, 2022.
- 377 [17] B. Kemp, A. Zwinderman, B. Tuk, H. Kamphuisen, and J. Obery. Analysis of a sleep-dependent
378 neuronal feedback loop: the slow-wave microcontinuity of the eeg. *IEEE Transactions on*
379 *Biomedical Engineering*, 47(9):1185–1194, 2000.
- 380 [18] T. Khandelwal, R. Das, and E. Chng. Sound event detection: A journey through dcase challenge
381 series. *APSIPA Transactions on Signal and Information Processing*, 13, 01 2024.
- 382 [19] G. Lavigne, P. Rompre, and J. Montplaisir. Sleep bruxism: validity of clinical research diagnostic
383 criteria in a controlled polysomnographic study. *Journal of dental research*, 75(1):546–552,
384 1996.
- 385 [20] W. T. McNicholas, D. Hansson, S. Schiza, and L. Grote. Sleep in chronic respiratory disease:
386 Copd and hypoventilation disorders. *European Respiratory Review*, 28(153), 2019.
- 387 [21] A. Mesaros, T. Heittola, and T. Virtanen. Metrics for polyphonic sound event detection. *Applied*
388 *Sciences*, 6:162, 05 2016.
- 389 [22] A. Mesaros, T. Heittola, T. Virtanen, and M. D. Plumbley. Sound event detection: A tutorial.
390 *IEEE Signal Processing Magazine*, 38(5):67–83, 2021.
- 391 [23] K. Mikkelsen, Y. Tabar, S. Kappel, C. Christensen, H. Toft, M. Hemmsen, M. Rank, M. Otto,
392 and P. Kidmose. Accurate whole-night sleep monitoring with dry-contact ear-eeg. *Scientific*
393 *Reports*, 9:16824, 11 2019.
- 394 [24] J.-K. Min, A. Doryab, J. Wiese, S. Amini, J. Zimmerman, and J. I. Hong. ‘Toss ’n’ turn:
395 smartphone as sleep and sleep quality detector. In *Proceedings of the SIGCHI Conference on*
396 *Human Factors in Computing Systems*, CHI ’14, page 477–486, New York, NY, USA, 2014.
397 Association for Computing Machinery.
- 398 [25] J. Monge-Álvarez, C. Hoyos-Barceló, P. Lesso, and P. Casaseca-de-la Higuera. Robust detection
399 of audio-cough events using local hu moments. *IEEE journal of biomedical and health*
400 *informatics*, 23(1):184–196, 2018.
- 401 [26] T. Nakamura, Y. Alqurashi, M. Morrell, and D. Mandic. Hearables: Automatic overnight sleep
402 monitoring with standardized in-ear eeg sensor. *IEEE Transactions on Biomedical Engineering*,
403 PP:1–1, 04 2019.
- 404 [27] R. Nandakumar, S. Gollakota, and N. Watson. Contactless sleep apnea detection on smartphones.
405 In *Proceedings of the 13th annual international conference on mobile systems, applications,*
406 *and services*, pages 45–57, 2015.
- 407 [28] A. Natsky, A. Vakulin, C. Coetzer, D. Mcevoy, R. Adams, and B. Kaambwa. Economic
408 evaluation of diagnostic sleep studies for obstructive sleep apnoea: a systematic review protocol.
409 *Systematic Reviews*, 10, 04 2021.
- 410 [29] J. Newell, O. Mairesse, P. Verbanck, and D. Neu. Is a one-night stay in the lab really enough
411 to conclude? First-night effect and night-to-night variability in polysomnographic recordings
412 among different clinical population samples. *Psychiatry Res*, 200(2-3):795–801, Dec 2012.
- 413 [30] A. Nobuyuki, N. Yasuhiro, T. Taiki, Y. Miyae, M. Kiyoko, and H. Terumasa. Trial of mea-
414 surement of sleep apnea syndrome with sound monitoring and spo2 at home. In *2009 11th*
415 *International Conference on e-Health Networking, Applications and Services (Healthcom)*,
416 pages 66–69. IEEE, 2009.

- 417 [31] C. O'Reilly, N. Gosselin, J. Carrier, and T. Nielsen. Montreal archive of sleep studies: an
418 open-access resource for instrument benchmarking and exploratory research. *Journal of sleep*
419 *research*, 23(6):628–635, 2014.
- 420 [32] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le. SpecAugment:
421 A simple data augmentation method for automatic speech recognition. In *Interspeech 2019*.
422 ISCA, sep 2019.
- 423 [33] Y. Ren, C. Wang, J. Yang, and Y. Chen. Fine-grained sleep monitoring: Hearing your breathing
424 with smartphones. In *2015 IEEE Conference on Computer Communications (INFOCOM)*,
425 pages 1194–1202. IEEE, 2015.
- 426 [34] T. Röddiger, C. Clarke, P. Breitling, T. Schneegans, H. Zhao, H. Gellersen, and M. Beigl. Sensing
427 with earables: A systematic literature review and taxonomy of phenomena. *Proceedings of the*
428 *ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 6(3):1–57, 2022.
- 429 [35] T. Röddiger, C. Dinse, and M. Beigl. Wearability and comfort of earables during sleep. In
430 *Proceedings of the 2021 ACM International Symposium on Wearable Computers, ISWC '21*,
431 page 150–152, New York, NY, USA, 2021. Association for Computing Machinery.
- 432 [36] M. J. Sateia. International classification of sleep disorders. *Chest*, 146(5):1387–1394, 2014.
- 433 [37] R. Serizel, N. Turpault, A. Shah, and J. Salamon. Sound event detection in synthetic domestic
434 environments. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech*
435 *and Signal Processing (ICASSP)*, pages 86–90, 2020.
- 436 [38] N. Shao, X. Li, and X. Li. Fine-tune the pretrained atst model for sound event detection. In
437 *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing*
438 *(ICASSP)*, pages 911–915, 2024.
- 439 [39] S. Shetty, V. Pitti, C. Satish Babu, G. Surendra Kumar, and B. Deepthi. Bruxism: a literature
440 review. *The Journal of Indian prosthodontic society*, 10:141–148, 2010.
- 441 [40] Y. Tabar, K. Mikkelsen, M. Rank, M. Hemmsen, M. Otto, and P. Kidmose. Ear-eeeg for sleep
442 assessment: a comparison with actigraphy and psg. *Sleep and Breathing*, 25:1–13, 09 2021.
- 443 [41] L. M. Trotti. Restless legs syndrome and sleep-related movement disorders. *Continuum:*
444 *Lifelong Learning in Neurology*, 23(4):1005–1016, 2017.
- 445 [42] N. Turpault, R. Serizel, J. Salamon, and A. P. Shah. Sound event detection in domestic
446 environments with weakly labeled data and soundscape synthesis. 4th Workshop on Detection
447 and Classification of Acoustic Scenes and Events (DCASE 2019), New York University, 2019.
- 448 [43] T. T. Um, F. M. J. Pfister, D. Pichler, S. Endo, M. Lang, S. Hirche, U. Fietzek, and D. Kulić. Data
449 augmentation of wearable sensor data for parkinson's disease monitoring using convolutional
450 neural networks. In *Proceedings of the 19th ACM International Conference on Multimodal*
451 *Interaction, ICMI '17*, page 216–220, New York, NY, USA, 2017. Association for Computing
452 Machinery.
- 453 [44] B. Venema, J. Schiefer, V. Blazek, N. Blanic, and S. Leonhardt. Evaluating innovative in-ear
454 pulse oximetry for unobtrusive cardiovascular and pulmonary monitoring during sleep. *IEEE*
455 *journal of translational engineering in health and medicine*, 1:2700208–2700208, 2013.
- 456 [45] Y. Wang, J. Ding, I. Chatterjee, F. Salemi Parizi, Y. Zhuang, Y. Yan, S. Patel, and Y. Shi. Faceori:
457 Tracking head position and orientation using ultrasonic ranging on earphones. In *Proceedings*
458 *of the 2022 CHI Conference on Human Factors in Computing Systems, CHI '22*, New York,
459 NY, USA, 2022. Association for Computing Machinery.

- 460 [46] Y. Wang, X. Zhang, J. M. Chakalasiya, X. Xu, Y. Jiang, Y. Li, S. Patel, and Y. Shi. Hearcough:
461 Enabling continuous cough event detection on edge computing hearables. *Methods*, 205:53–62,
462 2022.
- 463 [47] Y. Wu, K. Chen, T. Zhang, Y. Hui, T. Berg-Kirkpatrick, and S. Dubnov. Large-scale contrastive
464 language-audio pretraining with feature fusion and keyword-to-caption augmentation. In
465 *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing*
466 (*ICASSP*), pages 1–5. IEEE, 2023.
- 467 [48] X. Xu, E. Nemati, K. Vatanparvar, V. Nathan, T. Ahmed, M. Rahman, D. Mccaffrey, J. Kuang,
468 and J. U. N. A. Gao. Listen2Cough : Leveraging End-to-End Deep Learning Cough Detection
469 Model to Enhance Lung Health Assessment Using Passively Sensed Audio. *Proceedings of the*
470 *ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 5(1):1–22, 2021.
- 471 [49] Z. Yang, P. H. Pathak, Y. Zeng, X. Liran, and P. Mohapatra. Vital sign and sleep monitoring
472 using millimeter wave. *ACM Transactions on Sensor Networks (TOSN)*, 13(2):1–32, 2017.
- 473 [50] F. Zhang, C. Wu, B. Wang, M. Wu, D. Bugos, H. Zhang, and K. R. Liu. Smars: Sleep monitoring
474 via ambient radio signals. *IEEE Transactions on Mobile Computing*, 20(1):217–231, 2019.
- 475 [51] X. Zhang, Y. Wang, Y. Han, C. Liang, I. Chatterjee, J. Tang, X. Yi, S. Patel, and Y. Shi. The
476 earsavas dataset: Enabling subject-aware vocal activity sensing on earables. *Proc. ACM Interact.*
477 *Mob. Wearable Ubiquitous Technol.*, 8(2), may 2024.
- 478 [52] X. Zhang, Y. Wang, J. Zhang, Y. Yang, S. Patel, and Y. Shi. Earcough: Enabling continuous
479 subject cough event detection on hearables. In *Extended Abstracts of the 2023 CHI Conference*
480 *on Human Factors in Computing Systems*, pages 1–6, 2023.

481 A Experiment Detail

482 A.1 Hardware

483 The detailed hardware implementation is described in Table 6. Temperature was recorded but not
484 used during annotation and benchmarking models.

Module	Hardware	Frequency	Configuration Detail
dual-microphones IMU	Bose QC 20 MPU6090	24 kHz ≈94 Hz	I2S protocol, 16-bit PDM data format I2C protocol, 7-channel including accelerometer, gyro- scope and thermometer
compute chip	ESP32-WROOM	240 MHz	data transfer at 100KB/s with SPI protocol and an external TF card, 3.7V 800 mAH battery

Table 6: Hardware Implementation Detail.

485 A.2 Anonymized Participant Information

486 We recruited participants for our study through a campus study recruitment platform that included
487 students, faculty, and their family members. Table 7 shows the self-reported sleep disorder information
488 for all participants along with whether any apnea events were detected by those who wore the PDX
489 device.

Pair #	P1 Reported Disorder	P1 Apnea Detected (Y/N)	P2 Reported Disorder
1	Snore/Bruxism/Somniloquy	Y	—
2	Snore	Y	—
3	Bruxism	N	—
4	Snore/Bruxism	Y	—
5	—	N	—
6	Snore	N	—
7	Snore/Bruxism	N	Snore
8	Snore	Y	Snore
9	Bruxism/Somniloquy	N	—
10	Bruxism/Somniloquy	N	Somniloquy
11	Snore	Y	—

Table 7: Self-reported Sleep and Detected Sleep Disorders Among Participants.

490 A.3 Protocol and Compensation

491 To familiarize study participants with the data collection hardware and protocol, they were given a
492 video tutorial along with the following set of instructions:

- 493 1. Wear the PDx device according to the user guide and video to ensure that the sensors are working
494 properly.
- 495 2. Put on the headphones and start recording. To synchronize the data, please set an alarm on a
496 networked mobile phone for the nearest whole hour (e.g., 9 PM, 10 PM) after the headphones
497 begin recording. When the alarm sounds, please announce the time loud enough so that it can be
498 heard by both sets of headphones.
- 499 3. You should wear the headphones throughout the night for three consecutive nights, ensuring at
500 least 6 hours of data recording each night.
- 501 4. You should charge the devices after each recording so that they are ready for the next night’s
502 experiment.
- 503 5. After waking up each morning, please fill out a sleep diary to indicate when you fell asleep,
504 when you woke up, and any instances when you woke up during the night. In addition, please
505 fill out a PSQI sleep quality questionnaire.

506 Given that the hourly minimum wage was \$10 USD where this research was conducted, participants
 507 were paid \$70 USD per night of sleep. We collected 62 nights of data in total, totaling \$4,340 USD
 508 across the entire study.

509 **A.4 Manual Annotation**

510 Figure 5 shows the Audacity UI annotators used to examine and label data. The interface included
 511 each participant’s dual-channel audio and six-channel IMU data. Annotators created a separate
 512 "annotation" track where they could set the start and stop times of different events along with their
 513 labels. Table 8 shows the total number of events identified within each participant’s data.

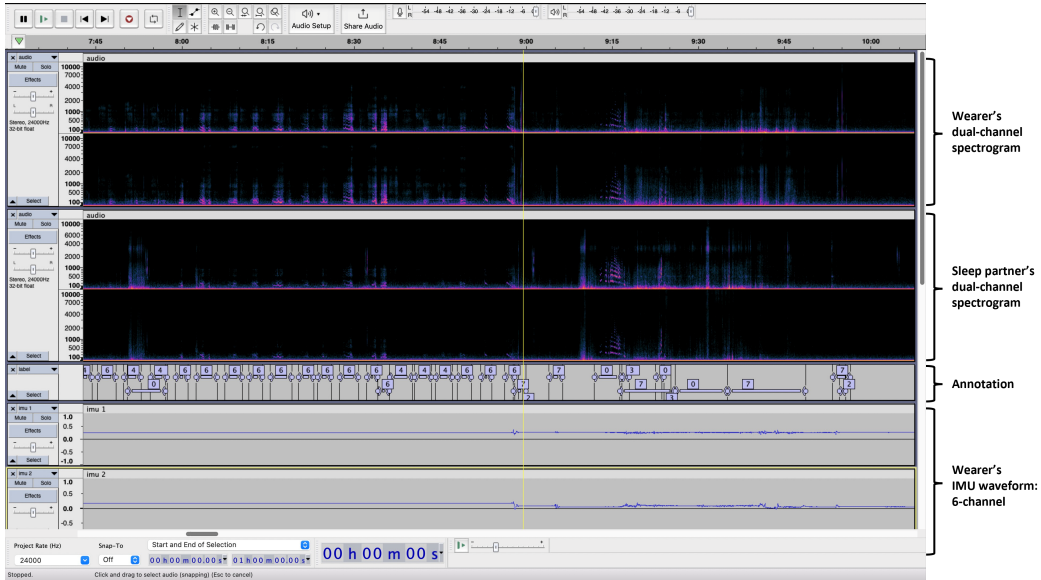
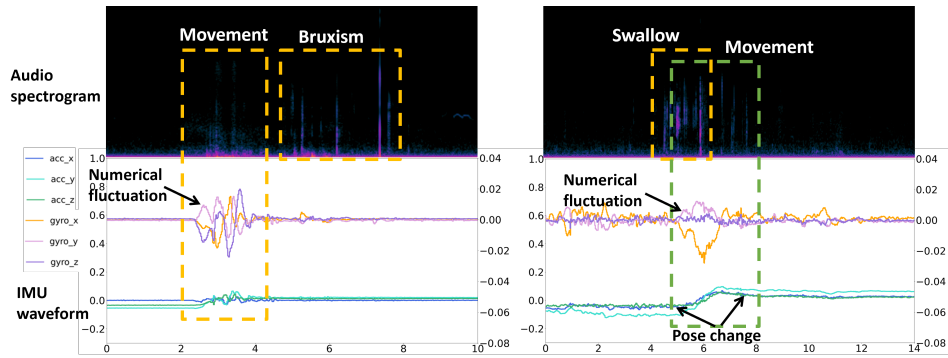


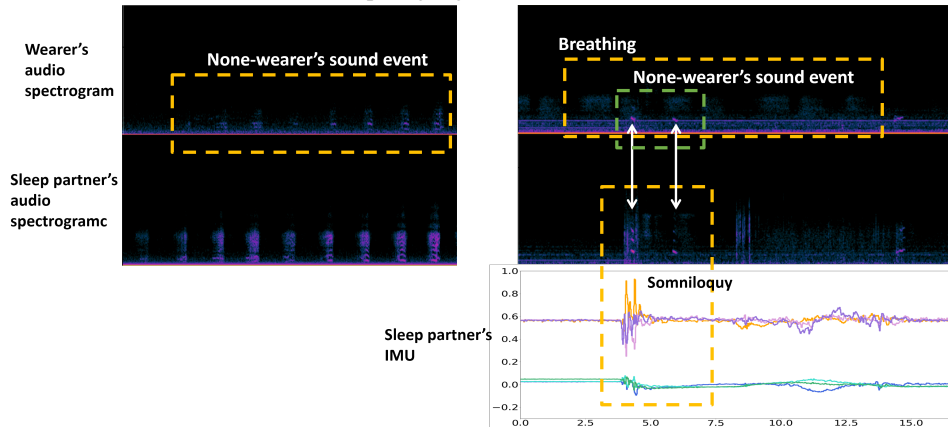
Figure 5: Annotation Software UI.

User	Label							
	Non-wearer's Sounds	Bruxism	Swallow	Somniloquy	Breathe	Cough	Snore	Movement
1-1	634	450	121	68	4565	32	13396	320
1-2	16720	399	214	86	2580	4	96	393
2-1	734	588	197	61	9629	6	3153	325
2-2	2400	118	116	100	3682	2	4	596
3-1	708	499	167	13	5684	2	328	266
3-2	882	245	90	10	1308	2	12	293
4-1	1059	478	319	6	8447	14	2419	217
4-2	5061	53	21	1	831	7	19	153
5-1	523	407	222	27	5589	1	18	438
5-2	614	599	131	26	9491	3	153	311
6-1	206	1050	102	36	12210	2	470	316
6-2	202	428	151	12	10248	3	484	226
7-1	7547	655	178	17	8280	3	2778	184
7-2	2600	240	183	36	8512	4	6126	246
8-1	43	81	26	6	2602	1	2245	47
8-2	2161	199	19	4	63	1	31	47
9-1	267	661	57	71	6954	3	328	148
9-2	2197	321	39	75	5708	0	11	177
10-1	316	519	85	93	8617	2	205	195
10-2	464	633	97	129	6160	0	26	320
11-1	1	43	43	6	578	0	1739	40
11-2	1847	37	10	1	1060	0	4	31
12-1	1058	658	62	30	9802	3	3178	162
12-2	4615	126	50	1	5071	0	7	139
Total	52859	9487	2700	915	137671	95	37230	5590

Table 8: Count of Sleep Events Per Participant.



(a) Comparing Signals Across Sensors.



(b) Comparing Audio Across Participants.

Figure 6: Aligned Data Annotation.

514 As shown in Figure 6, comparing data across wearers and non-wearers helped annotators identify the
 515 sources of different sounds. While loud breathing, snoring, or body movement by one individual can
 516 be captured by the earbuds of another in the same room, annotators were able to attribute the event's
 517 origin to the device that recorded the higher audio intensity. In some cases, the movement data also
 518 facilitated annotation. For example, the audio spectrograms of swallowing and body movement are
 519 quite similar and often occur simultaneously. However, the latter typically induces large fluctuations
 520 in IMU data due to changes in posture while swallowing is more subtle.

521 B Benchmark: Wearer Event Identification

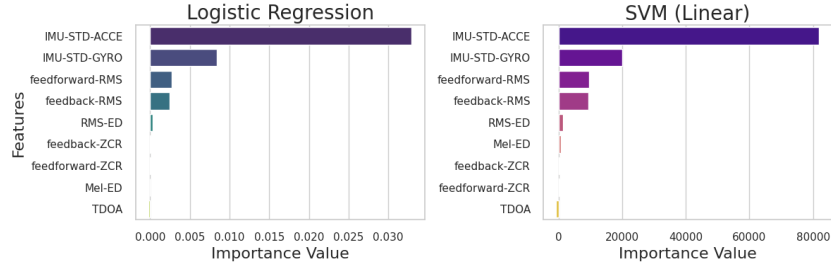
522 B.1 Supplementary Results

523 Figure 7 shows the importance of the audio and motion features we used in the lightweight and
 524 complex models for this task.

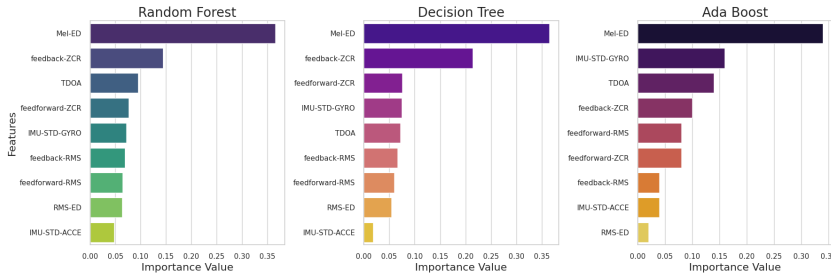
525 C Benchmark: Wearer-Aware Sleep Sound Event Classification

526 C.1 Model Architectures

527 All of the models were built using PyTorch 1.13.0 and trained on an NVIDIA GeForce RTX 4090
 528 GPU. We followed each paper as closely as possible and leveraged the accompanying code when
 529 available, using 5-second-long data inputs. The following describes the specific implementations we
 530 employed in this benchmark:



(a) Lightweight Models.



(b) Complex Models.

Figure 7: Feature Importances for the Wearer Event Identification Models.

531 **SleepTSM [11]:** We processed audio segments into log-scaled Mel filter bank features using a 0.1-s
 532 Hanning window with a 2-ms stride, yielding a 128×1501 input for each channel. Since SleepTSM
 533 is not open-sourced, we re-implemented it as described in the paper. The model was trained for 50
 534 epochs using the Adam optimizer with a constant learning rate of $1e-4$, a batch size of 32×4 , and a
 535 warmup ratio of 0.1.

536 **EarVAS [51]:** We processed audio segments into log-scaled Mel filter bank features using a 25-ms
 537 Hanning window with a 10-ms stride, yielding a $2 \times 498 \times 128$ input across both channels. As a
 538 requirement of the model’s EfficientNet-B0 sub-module, we implemented zero-padding on the filter-
 539 bank features along the time axis to obtain $2 \times 512 \times 128$ (channel \times time \times frequency) feature vectors.
 540 The motion input of EarVAS is raw 6-axis IMU data without any pre-processing. For the model
 541 itself, we leveraged the open-source code available at <https://github.com/thuhci/EarSAVAS>,
 542 which is shared under the MIT License found at [https://github.com/thuhci/EarSAVAS?tab=](https://github.com/thuhci/EarSAVAS?tab=MIT-1-ov-file)
 543 [MIT-1-ov-file](https://github.com/thuhci/EarSAVAS?tab=MIT-1-ov-file). The model used SpecAugment [32] for data augmentation and was trained for 30
 544 epochs using the Adam optimizer with a constant learning rate of $1e-4$ and a batch size of 128.

545 **Wav2Vec2.0 [5]:** We processed audio segments using the authors’ bespoke feature extractor. We
 546 then finetuned their open-sourced base model found at [https://huggingface.co/facebook/](https://huggingface.co/facebook/wav2vec2-base)
 547 [wav2vec2-base](https://huggingface.co/facebook/wav2vec2-base), which is shared under the Apache v2.0 License found at [https://huggingface.](https://huggingface.co/datasets/choosealicense/licenses/blob/main/markdown/apache-2.0.md)
 548 [co/datasets/choosealicense/licenses/blob/main/markdown/apache-2.0.md](https://huggingface.co/datasets/choosealicense/licenses/blob/main/markdown/apache-2.0.md). The
 549 model was trained for 10 epochs using the Adam optimizer with a constant learning rate of $3e-5$,
 550 a batch size of 32×4 , and a warmup ratio of 0.1.

551 **BEATs [10]:** We processed audio segments into log-scaled Mel filter bank features using a 25-ms
 552 Hanning window with a 10-ms stride, yielding 498×128 input for each channel. We leveraged the
 553 open-source model available at [https://github.com/thuhci/EarSAVAS/tree/main/BEATs_](https://github.com/thuhci/EarSAVAS/tree/main/BEATs_on_EarSAVAS)
 554 [on_EarSAVAS](https://github.com/thuhci/EarSAVAS/tree/main/BEATs_on_EarSAVAS), which is licensed under the MIT License found at [https://github.com/thuhci/](https://github.com/thuhci/EarSAVAS?tab=MIT-1-ov-file)
 555 [EarSAVAS?tab=MIT-1-ov-file](https://github.com/thuhci/EarSAVAS?tab=MIT-1-ov-file). The model was trained for 30 epochs using the Adam optimizer
 556 with a constant learning rate of $1e-4$ and a batch size of 32.

557 **CLAP [47]:** We performed zero-shot classification on the open-source model avail-
 558 able at https://huggingface.co/laion/larger_clap_general under the Apache v2.0
 559 License found at <https://huggingface.co/datasets/choosealicense/licenses/blob/main/markdown/apache-2.0.md>.

561 **C.2 Supplementary Results**

562 Figure 8 shows the confusion matrices associated with all of the models that were trained in this
 563 benchmark. The rest of this subsection describes the notable trends.

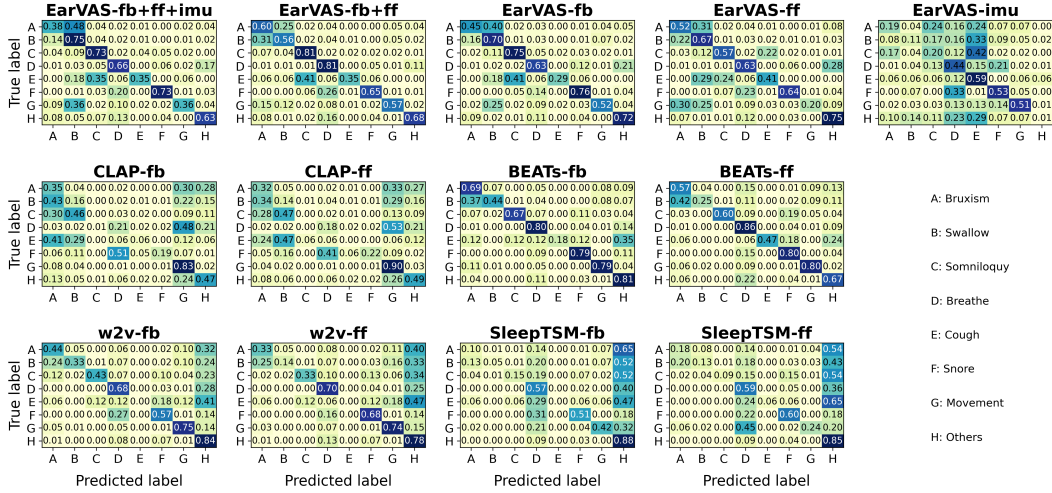


Figure 8: Confusion Matrices of the Wearer-aware Sleep Sound Event Classification Models.

564 **Challenges in class imbalance:** Coughing was the most underrepresented sleep event category in
 565 our dataset. Many of the models struggled to correctly identify these events, often confusing them
 566 with somnoliquy or external sounds. This underscores the necessity of accounting for dataset class
 567 imbalance in the wearer-aware sleep sound event classification. Although the class imbalance represents
 568 the real distribution of events in sleep scenarios, we have constructed a balanced dataset based on
 569 DreamCatcher to investigate whether models can achieve improved classification category-balanced
 570 conditions. The details are discussed in Appendix C.3.

571 **Challenges in jointly optimizing wearer event identification and sleep sound classification:**
 572 While traditional machine learning methods with feature engineering enable effective wearer event
 573 identification, wearer-aware sleep sound classification is a more challenging task because it also
 574 requires a model to perform sound event classification. While SleepTSM and CLAP have both been
 575 successfully used to discriminate between different sound events, introducing an additional class for
 576 background sounds and noises from the non-wearer resulted in poor accuracy. This underscores the
 577 necessity of carefully designing joint optimization methods specifically tailored for this task.

578 **Utility of the feedback microphone audio:** In Table 1, we found that benchmark models using
 579 only the feedback microphone audio exhibited a slightly improved performance compared to those
 580 that only used feedforward microphone audio. Upon further examination of the confusion matrices in
 581 Figure 8, we observe that utilizing feedback microphone audio can more effectively detect subtle
 582 swelling events that likely have lower intensity compared to other categories. This result highlights
 583 an important affordance of active noise-cancelling earbuds.

584 **Potential benefits of sensor fusion:** Comparing the performance of the EarVAS variants with
 585 different input feature modalities, we observed that leveraging both audio channels yielded higher

586 accuracy compared to either channel alone. Integrating motion data enhanced the model’s ability to
 587 discriminate wearer events from non-wearer events, which aligns with the conclusions drawn in the
 588 previous benchmark. However, the classification accuracy for categories like body movement and
 589 somniloquy decreased with this inclusion. With the proposed multi-modal DreamCatcher dataset, we
 590 encourage researchers to explore efficient sensor fusion methods for wearer-aware sleep sound event
 591 classification.

592 C.3 Wearer-Aware Sleep Sound Event Classification on Balanced Dataset

593 To investigate the influence on performance caused by the class imbalanced in wearer-aware sleep
 594 sound event classification, we constructed a balanced dataset based on DreamCatcher. Additionally, a
 595 comparison of the results from training Wav2Vec2.0 on both balanced and imbalanced datasets is
 596 conducted. The following describes the construction methods of the balanced dataset and notable
 597 points we found among the comparison.

598 **Construction of balanced training dataset based on DreamCatcher:** The categories and cor-
 599 responding number of samples in the imbalanced training dataset of DreamCatcher are as follows:
 600 Bruxism 6055, Swallow 1669, Somniloquy 612, Breathe 87898, Cough 61, Snore 23626, Movement
 601 3957. We set 10,000 as the target sample numbers for each category. For categories with more than
 602 10,000 samples, we randomly select 10,000 samples in each epoch. We also augment the categories
 603 with fewer than this number to 10,000.

604 The augmentation methods conducted to the rare events are detailed as below: in terms of audio
 605 augmentation, we referred to the augmentation methods in Audiomentations ⁸, which include 1)
 606 gain adjustment ($\times 0.5$ to $\times 2$), 2) time shift (-0.15 to 0.15 seconds), 3) pitch shift ($\times 0.5$ to $\times 2$), 4)
 607 speed adjustment ($\times 0.5$ to $\times 2$), and 5) random masking by making 0–10% of random points zero.
 608 As for motion data, we implemented augmentation according to the methods proposed by Terry et
 609 al. [43], including jittering, scaling, magnitude-warping, time-warping, rotations among three-axis
 610 accelerometer and three-axis gyroscope respectively, and permutation.

611 **Evaluation of Wav2Vec2.0 on balanced and imbalanced dataset:** We trained Wav2Vec2.0 on
 612 the balanced training dataset described above with the same settings as the Wav2Vec2.0 training on
 613 the imbalanced DreamCatcher. The two models are evaluated on the same testing dataset of the raw
 614 DreamCatcher. We present a comparative results of Wav2Vec2.0 before (w2v-before-bal) and after
 615 the balancing of training dataset (w2v-after-bal) in Table 9 and Figure 9.

Method	input channel	Evaluation Metrics (%)			
		Acc.	Macro-AUC	Macro-F1	MCC
w2v-before-bal	feedback	75.45	91.60	48.36	56.72
w2v-after-bal	feedback	70.70	89.93	46.04	48.76

Table 9: Comparative performance of Wav2Vec2.0 training on balanced and imbalanced dataset.

616 According to the results in the figure 9, we find that balanced dataset enhances the model’s ability
 617 to learn the patterns of rare events, thereby improving the recall rate. This is particularly notable
 618 for the ‘Cough’ category. However, data balancing also led to an overall decline in performance as
 619 shown in the table 9. This also highlights the diversity of samples from downsampling categories,
 620 making DreamCatcher a valuable dataset for sleep monitoring. We envision the generation of rare
 621 sleep events, which must aligns with the patterns of human sleep, will be a highly valuable area for
 622 future research.

⁸<https://github.com/iver56/audiomentations>

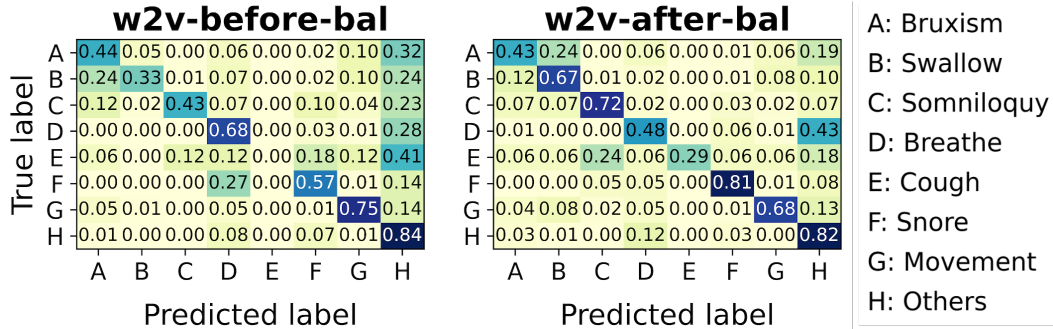


Figure 9: Confusion Matrices of Wav2Vec2.0 training on balanced and imbalanced dataset.

623 D Benchmark: Wearer-Aware Sleep Sound Event Detection

624 D.1 Model Architectures

625 Both models were built using PyTorch 1.13.0 and trained on an NVIDIA GeForce RTX 4090 GPU.
 626 We followed each paper as closely as possible and leveraged the accompanying code that was made
 627 available in both cases. Following the data format standards from the DCASE Challenges [18], we
 628 segmented the audio into 10-second clips and then generated the corresponding event onsets, offsets,
 629 and categories as labels.

630 **SEDNet [1]:** The original SEDNet leveraged audio sampled at 44.1 kHz, but we scaled its
 631 pre-processing steps to account for the fact that our audio was sampled at 16 kHz. Specif-
 632 ically, we processed audio segments by extracting log-scaled Mel-band energies using 40
 633 bands within a 2048-point Hamming window and a 1024-point stride, yielding a $2 \times 155 \times 40$
 634 (channel \times time \times frequency) feature vectors. We leveraged the open-sourced model found at <https://github.com/sharathadavanne/sed-crnn>, which is shared under the license found at <https://github.com/sharathadavanne/sed-crnn?tab=License-1-ov-file#readme>. The model
 635 was trained using the Adam optimizer with a batch size of 32. Training was stopped early if the
 636 Macro-F1 did not improve for 50 epochs.

639 **ATST-SED [38]:** For the model’s CNN module, we processed audio segments by extracting 128
 640 log-scaled Mel features from frames with a 128-ms length and a 16-ms stride. For the model’s
 641 ATST-Frame module, we converted the audio segments into log-Mel spectrograms using a 64-ms
 642 Hamming window with a 10-ms stride. The resulting spectrogram comprised 64 Mel-frequency
 643 bins spanning a frequency range from 60 Hz to 7800 Hz. We leveraged the open-source model
 644 found at <https://github.com/Audio-WestlakeU/ATST-SED>, which is shared under the MIT
 645 License found at <https://github.com/Audio-WestlakeU/ATST-SED?tab=MIT-1-ov-file>.
 646 The model was trained using the Adam optimizer with a batch size of 24. We trained the first stage of
 647 ATST-SED for 200 epochs, during which the pretrained ATST-Frame module was frozen while the
 648 remaining parts were trained. Due to significant performance degradation observed during the second
 649 stage of training, we did not utilize it in this benchmark.

650 D.2 Evaluation Metrics

651 We evaluated the models in this benchmark using collar-based metrics [21] that compare the onset and
 652 offset of a predicted and target event. An offset is often used to account for differences in prediction
 653 resolution. Based on the empirical standard proposed by Serizel et al. [37] for sound event detection,
 654 we used a 200 ms tolerance to compare onset timestamps, and we used the maximum of 200 ms and
 655 20% of the duration of the sound event to compare offset timestamps. With these considerations in
 656 mind, we calculated the following performance metrics:

657 **Macro-F1.** The macro-averaged F1 score is the arithmetic mean of F1 scores in a multi-class model.
658 In this context, true positives are defined as events in the system output that have a temporal position
659 overlapping with the temporal position of an event with the same label in the ground truth. False
660 positives are defined as events in the system output that have no correspondence to an event with
661 the same label in the ground truth within the allowed tolerance, while false negatives are defined as
662 events in the ground truth that have no correspondence to an event with the same label in the system
663 output within the allowed tolerance. These metrics are computed on a per-class basis to produce
664 class-specific F1 scores that are combined to calculate Macro-F1.

665 **Error Rate.** Error rate is calculated as the total number of substitutions, deletions, and insertions
666 divided by the total number of events in the ground truth. Substitutions are events in the system
667 output with a correct temporal position but an incorrect class label, insertions are extraneous events
668 in the system output, and deletions are events in ground truth not included in the system output.

669 **D.3 Supplementary Results**

670 Table 10 examines the class-wise performance of the models that were trained in this benchmark.
671 The rest of this subsection describes the notable trends

672 **Challenges in class imbalance:** Similar to the supplementary results for the previous benchmark,
673 class imbalance also had an impact on sound event detection. Both models performed poorly on the
674 severely undersampled cough class, to the point that they were never able to detect such events in
675 the dataset. However, the data imbalance issue extends beyond different categories of events. It also
676 includes variation in the duration of individual samples and imbalances between active and inactive
677 frames [16], which caused notable challenges for SEDNet specifically. Since our dataset reflects the
678 natural distribution of sleep events in non-restrictive environments, we encourage the development of
679 targeted solutions to address these challenges.

680 **Utility of the feedback microphone audio:** In Table 5, we found that benchmark models using
681 only the feedback microphone audio exhibited better performance compared to those that only used
682 feedforward microphone audio. Upon further examination of the class-wise accuracies in Table 10,
683 we observed that the high signal-to-noise ratio of the feedback microphone audio enables the model to
684 detect subtle sounds. This was particularly evident with ATST-SED, which saw a marked performance
685 boost in classes like bruxism.

686 **Potential benefits of sensor fusion.** Since SEDNet supports multi-channel audio, we used it to
687 examine the utility of dual-channel audio fusion methods. However, we found that this approach
688 actually resulted in performance degradation. This may be because SEDNet was designed to
689 utilize a multi-channel microphone array to localize sound sources, whereas the microphones in
690 this application are in much closer proximity to one another and the sound source. The EarVAS
691 model [51] for subject-aware vocal activity classification demonstrates that models intentionally
692 trained for earbud hardware can overcome this issue. Furthermore, we did not investigate the utility
693 of using IMU data for event detection, but given its utility in wearer event identification, it may be
694 fruitful to pursue this direction further in a multi-modal architecture.

Event label	F1	Pre.	Rec.	ER.	Del.	Ins.
Bruxism	12.7%	30.4%	8.0%	1.10	0.92	0.18
Swallow	4.1%	33.3%	2.2%	1.02	0.98	0.04
Somniloquy	5.0%	14.3%	3.0%	1.15	0.97	0.18
Breathe	56.2%	56.6%	55.7%	0.87	0.44	0.43
Cough	0.0%	0.0%	0.0%	1.00	1.00	0.00
Snore	52.1%	55.6%	49.0%	0.90	0.51	0.39
Movement	38.7%	41.2%	36.4%	1.15	0.64	0.52

(a) ATST-SED With Feedforward Microphone Audio.

Event label	F1	Pre.	Rec.	ER.	Del.	Ins.
Bruxism	21.6%	30.8%	16.7%	1.21	0.83	0.37
Swallow	6.9%	16.4%	4.4%	1.18	0.96	0.22
Somniloquy	2.8%	20.0%	1.5%	1.05	0.98	0.06
Breathe	57.6%	56.7%	58.5%	0.86	0.41	0.45
Cough	0.0%	0.0%	0.0%	1.00	1.00	0.00
Snore	47.4%	54.3%	42.0%	0.93	0.58	0.35
Movement	36.8%	42.1%	32.6%	1.12	0.67	0.45

(b) ATST-SED With Feedback Microphone Audio.

Event label	F1	Pre.	Rec.	ER.	Del.	Ins.
Bruxism	0.0%	0.0%	0.0%	1.01	1.00	0.01
Swallow	0.0%	0.0%	0.0%	1.01	1.00	0.01
Somniloquy	0.0%	0.0%	0.0%	1.02	1.00	0.02
Breathe	50.8%	50.4%	51.2%	0.99	0.49	0.50
Cough	0.0%	0.0%	0.0%	1.11	1.00	0.11
Snore	44.2%	48.4%	40.6%	1.03	0.59	0.43
Movement	3.2%	3.3%	3.0%	1.86	0.97	0.89

(c) SEDNet With Dual-channel Audio.

Event label	F1	Pre.	Rec.	ER.	Del.	Ins.
Bruxism	0.0%	0.0%	0.0%	1.00	1.00	0.00
Swallow	0.0%	0.0%	0.0%	1.00	1.00	0.00
Somniloquy	0.0%	0.0%	0.0%	1.00	1.00	0.00
Breathe	49.5%	48.9%	50.1%	1.02	0.50	0.52
Cough	0.0%	0.0%	0.0%	1.00	1.00	0.00
Snore	44.6%	48.1%	41.6%	1.03	0.58	0.45
Movement	31.7%	31.5%	31.9%	1.37	0.68	0.69

(d) SEDNet With Feedforward Microphone Audio.

Event label	F1	Pre.	Rec.	ER.	Del.	Ins.
Bruxism	0.0%	0.0%	0.0%	1.00	1.00	0.00
Swallow	0.0%	0.0%	0.0%	1.00	1.00	0.00
Somniloquy	0.0%	0.0%	0.0%	1.00	1.00	0.00
Breathe	52.7%	52.9%	52.6%	0.94	0.47	0.47
Cough	0.0%	0.0%	0.0%	1.00	1.00	0.00
Snore	47.7%	48.9%	46.5%	1.02	0.54	0.49
Movement	31.6%	35.4%	28.5%	1.24	0.71	0.52

(e) SEDNet With Feedback Microphone Audio.

F1: F1 Score; Pre.: Precision, Rec.: Recall; ER.: Error Rate; Del.: Deletion Rate; Ins.: Insertion Rate

Table 10: Class-wise Results of the Wearer-aware Sleep Sound Event Detection Models.

695 Checklist

696 1. For all authors...

- 697 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
698 contributions and scope? [Yes]
- 699 (b) Did you describe the limitations of your work? [Yes]
- 700 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 701 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
702 them? [Yes]
- 703 2. If you are including theoretical results...
- 704 (a) Did you state the full set of assumptions of all theoretical results? [N/A] Did not
705 include theoretical results.
- 706 (b) Did you include complete proofs of all theoretical results? [N/A]
- 707 3. If you ran experiments (e.g. for benchmarks)...
- 708 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
709 mental results (either in the supplemental material or as a URL)? [Yes]
- 710 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
711 were chosen)? [Yes] See Section 3.2, C.1, D.1
- 712 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
713 iments multiple times)? [No] We do not have enough compute resources to repeat
714 experiments for error bars.
- 715 (d) Did you include the total amount of compute and the type of resources used (e.g., type
716 of GPUs, internal cluster, or cloud provider)? [Yes] See Section C.1, D.1
- 717 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 718 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 719 (b) Did you mention the license of the assets? [Yes]
- 720 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 721 (d) Did you discuss whether and how consent was obtained from people whose data you're
722 using/curating? [Yes] See Section 3.3
- 723 (e) Did you discuss whether the data you are using/curating contains personally identifiable
724 information or offensive content? [Yes] See Section 3.3
- 725 5. If you used crowdsourcing or conducted research with human subjects...
- 726 (a) Did you include the full text of instructions given to participants and screenshots, if
727 applicable? [Yes] See Section A.3
- 728 (b) Did you describe any potential participant risks, with links to Institutional Review
729 Board (IRB) approvals, if applicable? [Yes] See Section 3.3
- 730 (c) Did you include the estimated hourly wage paid to participants and the total amount
731 spent on participant compensation? [Yes] See Section A.3