

<https://doi.org/10.1038/s41746-024-01139-z>

# Towards automatic home-based sleep apnea estimation using deep learning



Gabriela Retamales<sup>1,4</sup>, Marino E. Gavidia<sup>1,4</sup>, Ben Bausch<sup>1</sup>, Arthur N. Montanari<sup>1,2</sup>,  
Andreas Husch<sup>1</sup> & Jorge Goncalves<sup>1,3</sup> ✉

Apnea and hypopnea are common sleep disorders characterized by the obstruction of the airways. Polysomnography (PSG) is a sleep study typically used to compute the Apnea-Hypopnea Index (AHI), the number of times a person has apnea or certain types of hypopnea per hour of sleep, and diagnose the severity of the sleep disorder. Early detection and treatment of apnea can significantly reduce morbidity and mortality. However, long-term PSG monitoring is unfeasible as it is costly and uncomfortable for patients. To address these issues, we propose a method, named DRIVEN, to estimate AHI at home from wearable devices and detect when apnea, hypopnea, and periods of wakefulness occur throughout the night. The method can therefore assist physicians in diagnosing the severity of apneas. Patients can wear a single sensor or a combination of sensors that can be easily measured at home: abdominal movement, thoracic movement, or pulse oximetry. For example, using only two sensors, DRIVEN correctly classifies 72.4% of all test patients into one of the four AHI classes, with 99.3% either correctly classified or placed one class away from the true one. This is a reasonable trade-off between the model's performance and the patient's comfort. We use publicly available data from three large sleep studies with a total of 14,370 recordings. DRIVEN consists of a combination of deep convolutional neural networks and a light-gradient-boost machine for classification. It can be implemented for automatic estimation of AHI in unsupervised long-term home monitoring systems, reducing costs to healthcare systems and improving patient care.

Apnea and hypopnea are common sleep-related breathing disorders that affect 6% to 17% of adults, reaching up to 49% in older populations<sup>1</sup>. They are usually marked by pauses in breathing or shallow breathing, respectively. This can happen, for example, when the soft tissue at the back of the throat collapses and closes at night<sup>2</sup>. Because blockage in the airways slows breathing, less oxygen is sent from the lungs to the heart and body. The CO<sub>2</sub> level in the blood then increases due to impaired breathing, leading to episodes of sudden awakening or choking during sleep<sup>3</sup>. During these events, nasal airflow can cease or reduce for a short time, and the brain and body will fight to keep breathing. Sleep apnea and hypopnea have been associated with an increased risk of heart disease, diabetes, chronic kidney disease, stroke, depression, and cognitive impairment<sup>2</sup>. As age and the likelihood of having sleep apnea are positively correlated<sup>4</sup>, an aging population is expected to add more pressure on healthcare systems.

Typically, the diagnosis and detection of sleep apnea and hypopnea are based on polysomnography (PSG) tests conducted in a sleep facility or at

home<sup>5</sup>. PSG requires the overnight recording and monitoring of several physiological signals, including electroencephalogram (EEG), electrocardiogram (ECG), electrooculogram, chin muscle activity, leg movements, respiratory effort, nasal airflow, and oxygen saturation (SpO<sub>2</sub>). This information is then used to compute the Apnea-Hypopnea Index (AHI), the number of times a person has apnea or certain types of hypopnea per hour of sleep, as defined by the American Academy of Sleep Medicine (AASM)<sup>6</sup>.

There are many limitations with PSG tests: they are time-consuming, expensive, and inconvenient<sup>7</sup>. As a result, it is expected that more than 85% of people with sleep apnea or hypopnea are not diagnosed<sup>4</sup>. Hence, the development of portable, easy-to-use, reliable, and affordable AHI estimation tools is crucial to improve patient care. Currently, there is only one method that estimates AHI<sup>8</sup>, using an oximetry signal. However, this method does not estimate when a patient is sleeping, which is essential to compute AHI. Instead, it uses the information of when patients are asleep from the labels provided in the dataset (ground truth). As we will show later,

<sup>1</sup>Luxembourg Centre for Systems Biomedicine, University of Luxembourg, L-4367, Belvaux, Luxembourg. <sup>2</sup>Department of Physics and Astronomy, Northwestern University, Evanston, IL 60208, USA. <sup>3</sup>Department of Plant Sciences, University of Cambridge, Cambridge CB2 3EA, UK. <sup>4</sup>These authors contributed equally: Gabriela Retamales, Marino E. Gavidia. ✉e-mail: [jmg77@cam.ac.uk](mailto:jmg77@cam.ac.uk)

oximetry signals on their own struggle to identify when patients are asleep, even when coupled with heart rate signals. This limitation can lead to poor AHI estimation. Finally, the available method does not detect events or segment oximetry signals, i.e., it does not identify the exact periods of time when apneas and hypopneas occur.

When focusing only on short time-series segments, there are many classification tools to decide whether an apnea event is present or not (see, for example, the comprehensive reviews in refs. 9–11). Most results are limited to obstructive sleep apnea (OSA), the most severe and common type of sleep apnea<sup>12</sup>. Such detection tools for OSA classification have been based on manual feature extraction, from time and/or frequency domains of one or more physiological signals. These algorithms require expert knowledge for the design and selection of features to be extracted, which may still not include the most relevant features (information) for classification from data<sup>13</sup>. The procedure can be hard, time consuming, and often subject to bias or lack of generalizability<sup>14,15</sup>. Deep learning, on the other hand, has shifted data modeling away from “expert-driven” feature engineering toward “data-driven” feature extraction<sup>8,11</sup> and is now widespread across many other biomedical applications<sup>16,17</sup>.

Based on deep learning, this paper proposes a data-driven method, named DRIVEN, to estimate AHI and segment physiological signals. We focus on devices that can be implemented in a home care environment and were available in the considered datasets, including abdominal and thoracic movement, oxygen saturation, and R-to-R interval data. Our study employs an extensive multicenter database containing 14,370 recordings of 10,752 patients belonging to multiple ethnicities and with an average age of 68.7 years. To estimate AHI, DRIVEN detects apnea and hypopnea events, and determines when patients are sleeping or awake. To our knowledge, it is the first standalone algorithm to estimate AHI. Using abdominal movement and oximetry sensors, DRIVEN correctly classifies 72.4% of all test patients into one of the four AHI classes, with 99.3% either correctly classified or placed one class away from the true one, a reasonable trade-off between the model’s performance and patient’s comfort. Cross-validation and out-of-the-domain testing show that DRIVEN is highly generalizable, achieving high performance on large external datasets not included during the training process. Given the high performance and small computational cost of DRIVEN, we expect that it could be implemented as part of a home AHI estimator system.

## Results

### Data description

Three different datasets were considered in this work: the Multi-Ethnic Study of Atherosclerosis (MESA)<sup>18</sup>, the Men Study of Osteoporotic Fracture (MrOS)<sup>19</sup>, and the Sleep Heart Health Research (SHHS)<sup>20</sup>. These datasets are publicly available from the National Sleep Research Resource for Sleep-Related Studies (NSRR)<sup>21</sup>. In total, 14,370 PSG recordings were considered in our study. Each PSG recording consists of data collected from several physiological sensors during a patient’s overnight sleep (lasting 8.78 h on average). The PSG recordings were collected during typical PSG evaluations in healthcare facilities.

This study focuses only on channels that can be easily implemented in a home environment and are accessible in all three datasets: pulse oximetry (SpO<sub>2</sub>), thoracic movement, and abdominal movement. Data collected from various healthcare facilities had different sampling frequencies. These were normalized by resampling all channels to the maximum sampling frequency (64 Hz). Hence, all signals per patient were standardized to the same length, which simplified the training of the neural networks. PSG recordings were excluded from the study if they met at least one of the following criteria: (1) poor-quality PSG recordings given by the presence of “unsure” or “noise” labels in more than a third of the total sleep time or the repetition of only one value through the whole signal; and (2) the absence of one or more channels from the five initial channels selected for the study. After applying the exclusion criteria, a total of 10,643 sleep recordings were obtained to develop and evaluate the proposed method. Table 1 summarizes the databases.

**Table 1 | Data description**

Characteristics	SHHS	MROS	MESA
Total PSG recordings	8444	3871	2055
Male	4458	3871	953
Female	3986	0	1102
Age	64.6 (39–90)	77.5 (67–90)	69.4 (54–94)
BMI	28.1 (18–50)	27.1 (16–47)	28.7 (16–56)
Sleep time [min]	364.34 (34–605)	352.19 (39–623)	359.72 (32–601)
AHI [per hour]	18.11 (0.0–161.8)	21.57 (0.0–106.0)	24.12 (0.0–111.9)
<i>Ethnicity</i>			
White	85.4%	89.8%	36.2%
African	8.2%	3.5%	27.8%
Hispanic	–	2.1%	23.9%
Asian	–	3.3%	12.2%
Other	6.4%	1.2%	–
Used PSG recordings	5288	3364	1991

Data are reported in mean and range (in parenthesis).

AHI Apnea-Hypopnea Index, BMI body mass index.

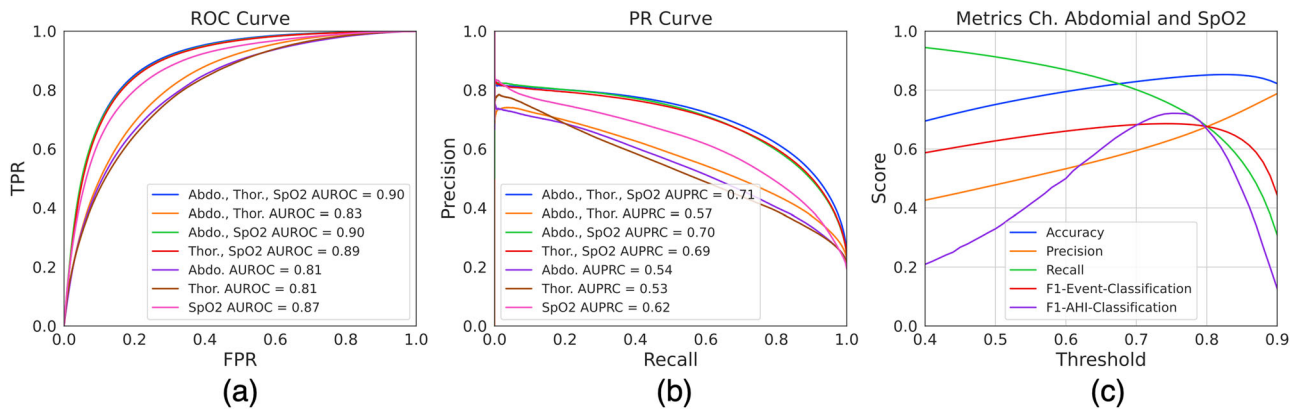
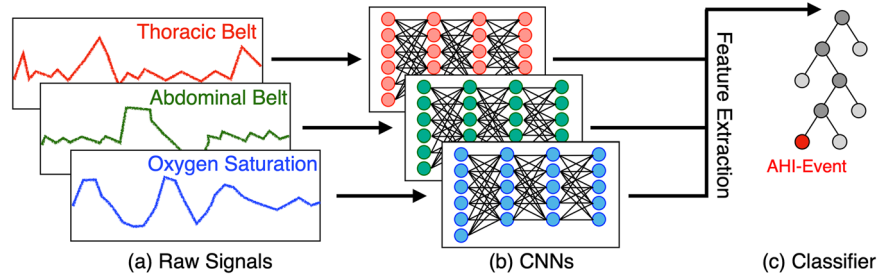
The data was labeled following the definition of AHI (ahi\_a0h3a) from the American Academy of Sleep Medicine<sup>6</sup>. An apnea is defined as a nasal flow reduction of more than 90%, while a hypopnea is defined as a nasal flow reduction of more than 30%, both for at least 10 s. There are several types of hypopneas, in which only the following two count towards AHI: hypopneas with a reduction of oxygen desaturation equal to 3% or more within 45 s (also labeled in the paper as hypopnea type 1) and hypopneas with an arousal within 5 s after the end of the event (hypopnea type 2)<sup>22</sup>. All other hypopneas (hypopnea type 3) are not considered in the calculation of AHI. In this paper, an *AHI event* refers to any of these events: apnea, hypopnea type 1, and hypopnea type 2. According to the AASM guidelines, patients are categorized according to their AHI: healthy (less than 5 events/h), mild (5–15 events/h), moderate (15–30 events/h), and severe (more than 30 events/h).

### DRIVEN: deep-learning model for the detection of AHI events

We developed a hybrid deep-learning model named DRIVEN to (1) classify “normal” or AHI events, (apneas, hypopneas type 1, and hypopneas type 2) from data, (2) estimate the total time of sleep, and (3) estimate the AHI. In the classification step, the inputs of DRIVEN are short segments of time-series data recorded during sleep. We explored window lengths of 10, 30, and 80 s, and selected 30 s windows as they outperformed the others at estimating AHI (see Methods for details). We focused on the following physiological signals: SpO<sub>2</sub>, thoracic movement, and abdominal movement. DRIVEN then learns whether a given input corresponds to a normal event or an AHI event. Figure 1 illustrates DRIVEN’s pipeline. For each available physiological signal, a distinct deep convolutional neural network (CNN) is trained to extract global features from the raw 1-dimensional signals. The extracted features from all CNNs are then combined on a light gradient-boost machine (LightGBM)<sup>23</sup> to perform the classification between normal and AHI events. This approach has proven to generalize better than just using one CNN to input all the channels (see subsection “End-to-end CNN” in Methods). Finally, for each patient, DRIVEN counts all detected AHI events to estimate the AHI. We evaluate the performance of different sensors, window sizes, and deep-learning models. See Methods for details on data pre-processing and model training, validation, and testing.

DRIVEN was trained and assessed on raw physiological signals. This avoids creating human-selected features, which can sometimes lead to bias. The data were separated by center to avoid data leakage and guarantee good performance on external datasets. SHHS was used to train and validate the model (both CNNs and LightGBM) while MESA and MROS were used as

**Fig. 1 | DRIVEN’s pipeline.** **a** Data are separated by channels and segmented into 30 s windows. **b** For each channel, a distinct trained deep CNN extracts features (outputs) from the raw signal (input). **c** The extracted features are concatenated and fed to a trained LightGBM that classifies the input data between normal and AHI events (apneas, hypopneas type 1, and hypopneas type 2).



**Fig. 2 | Performance of DRIVEN on the classification of AHI events.** **a** Receiver-operator characteristic and **(b)** precision-recall curves. Note the overlapping curves for thoracic and abdominal sensors. **c** Threshold-dependent performance metrics of DRIVEN when using two input channels (abdominal movement and SpO<sub>2</sub>). The

performance results are shown for all patients in the test dataset. The accuracy, precision, recall, and F1-Event-Classification are metrics on the classification of individual events. The F1-AHI-Classification measures the F1-Score of predicting the AHI severity category (healthy, mild, moderate, severe) on whole sleep studies.

test sets to independently evaluate DRIVEN on unseen out-of-domain data from other sources. For each patient, sleep intervals were sampled every 30 s with a 15 s overlap and labeled as normal or AHI event, yielding a total of more than 15 million samples. The training and validation data (SHHS database) were originally imbalanced with a ratio of 6.4 normal events per AHI event. To overcome the bias and inaccuracy associated with classification using imbalanced data between two classes, the number of normal events was under sampled. Balancing the number of positive and negative samples is a common strategy followed in the literature<sup>24</sup>.

**Performance of DRIVEN on sleep/awake detection**

The AHI score is defined as the quotient between the total number of AHI events recorded overnight and the total time that patient slept (see subsection “AHI and severity class estimation from positive events” in Methods). Hence, to get an accurate estimation of this number, DRIVEN also needs to estimate the time the patient slept through the night. We trained the same architecture proposed in Fig. 1, using window sizes of 30 s, with the annotated labels from the dataset for awake or sleep. The training, validation, and testing methodology were similar as before. For the combination of abdominal and SpO<sub>2</sub> sensors, we obtained an area under the receiver operating characteristic curve (AUROC) of 0.87 and an area under the precision-recall curve (AUPRC) of 0.90 in the test dataset. We also observe that adding RRI to any combination of sensors has little effect on the performance of the algorithm and that pulse oximetry (SpO<sub>2</sub>) alone is a poor sleep state predictor. The complete results of the sleep/wake classification are reported in Supplementary Fig. 1 and Supplementary Table 1.

**Performance of DRIVEN on classification of AHI events**

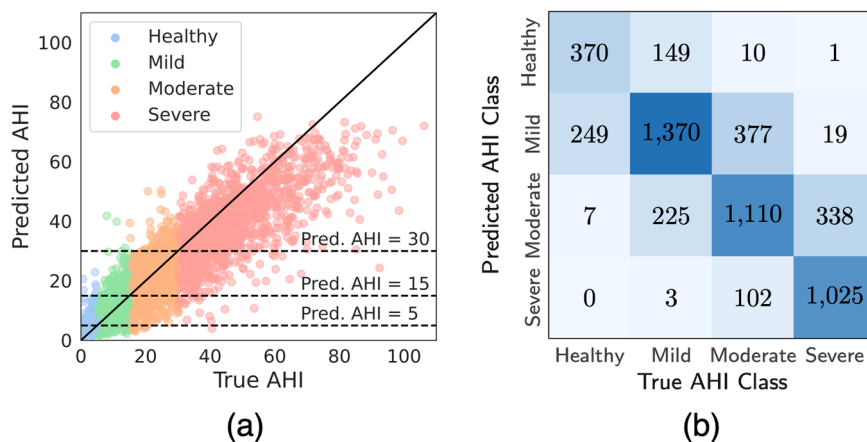
We now evaluate the performance of DRIVEN in the classification task between normal and AHI events. Results for the complete test set (MESA and MROS databases) are shown in Fig. 2a,b using either a single channel or a combinations of channels (when the sleep/awake states are also predicted

by DRIVEN). AUROC and AUPRC show that models using thoracic and abdominal movement sensors each feature very similar performance, even when these two sensors are used in combination. This indicates very high mutual information between these sensors. On the other hand, adding SpO<sub>2</sub> significantly improves the performance, indicating additional information carried by this signal. Note that the model using the SpO<sub>2</sub> channel alone struggles to classify AHI events due to its poor performance to estimate sleep/awake segments (Supplementary Fig. 1 shows a relatively low AUROC and AUPR of 0.72 and 0.76, respectively, on the sleep/awake detection when using only SpO<sub>2</sub>). Figure 2c shows performance metrics as a function of the threshold when using two sensors: abdominal movement and SpO<sub>2</sub>. Both accuracy and F1-Event-Classification curves are relatively flat for a wide range of thresholds showing that the model is robust to the choice of threshold (which was decided on the validation results). Note that estimating hypopneas type 2 requires detecting short-term arousals, which are hard to estimate from the three available sensors. In contrast, apneas and hypopneas type 1 are much easier to detect (Supplementary Figs. 2 and 11).

**Performance of DRIVEN on AHI estimation**

This section evaluates the performance of DRIVEN at estimating the AHI for each patient on the test set. It consists of first estimating the total sleep time and then counting all detected AHI events throughout the night for each patient. Using both abdominal movement and SpO<sub>2</sub> sensors, Fig. 3a shows the real versus predicted AHI while Fig. 3b contains the confusion matrix of the four severity categories. In total, 72.4% of patients are correctly classified and 99.3% were either correctly classified or were placed one class away from the true one. No healthy patient was classified as severe. Among the 5355 test patients, only one severe individual was classified as healthy. This occurred because the patient had a short sleep duration of 3.5 hours with frequent apneic bursts. DRIVEN exhibits consistent performance across all AHI severity levels, effectively identifying sleep-disordered breathing events in all patient groups. The model provides a good trade-

**Fig. 3 | Performance of DRIVEN on the estimation of AHI.** **a** Real versus predicted AHI divided by the four AHI severity groups. **b** Confusion matrix. The performance is evaluated on the test data considering a threshold of 0.79 and a combination of two signals (abdominal movement and SpO<sub>2</sub>).



off between accuracy and patient’s sleep comfort, when compared to a full polysomnography. Results for other combination of sensors are reported in Supplementary Figs. 3–8.

To investigate whether some physiological signals may contribute more than others in estimating AHI, we calculated F1 scores for different combinations of sensors (Table 2; further details can be found in Supplementary Figs. 3–8). Similar to the classification results above, thoracic and abdominal movement sensors each have similar performance, while adding SpO<sub>2</sub> significantly improves the performance. The combination of abdominal movement and SpO<sub>2</sub> sensors results in one of the highest performance and using only two sensors. We also compared the AHI classification F1 scores when AHI is estimated (1) without sleep information, (2) with ground truth on the sleep/awake states, and (3) with the predicted sleep/awake states (Supplementary Table 2). Using abdominal and SpO<sub>2</sub> sensors, the F1 performance drops around 4% from the ground-truth sleeping information to the predicted sleep. For SpO<sub>2</sub> alone, however, the drop is considerably larger at 12%, as this sensor alone cannot estimate sleep states as well.

**Automatic detection of AHI events: segmentation of whole sleep studies**

Using DRIVEN’s classification of AHI events, we can analyze whole sleep studies to identify apnea and certain hypopneas (types 1 and 2). Additionally, with DRIVEN’s sleep state detection, we can also identify when the patient is asleep and when the patient is awake. For each patient, recordings were sampled every 15 s with 30 s windows. Each window was labeled as awake or sleep according to the sleep prediction classification. Then, those segments classified as sleep are further categorized into normal or AHI events. Figure 4 illustrates the method on a random patient, revealing the periods of sleep with the highest AHI event episodes, compared with the ground truth, as well as awake segments.

**Discussion**

This paper introduced a method, DRIVEN, to (1) classify AHI events and sleep states, (2) estimate AHI, and (3) segment whole sleep studies, using a

**Table 2 | F1 scores for different combinations of sensors**

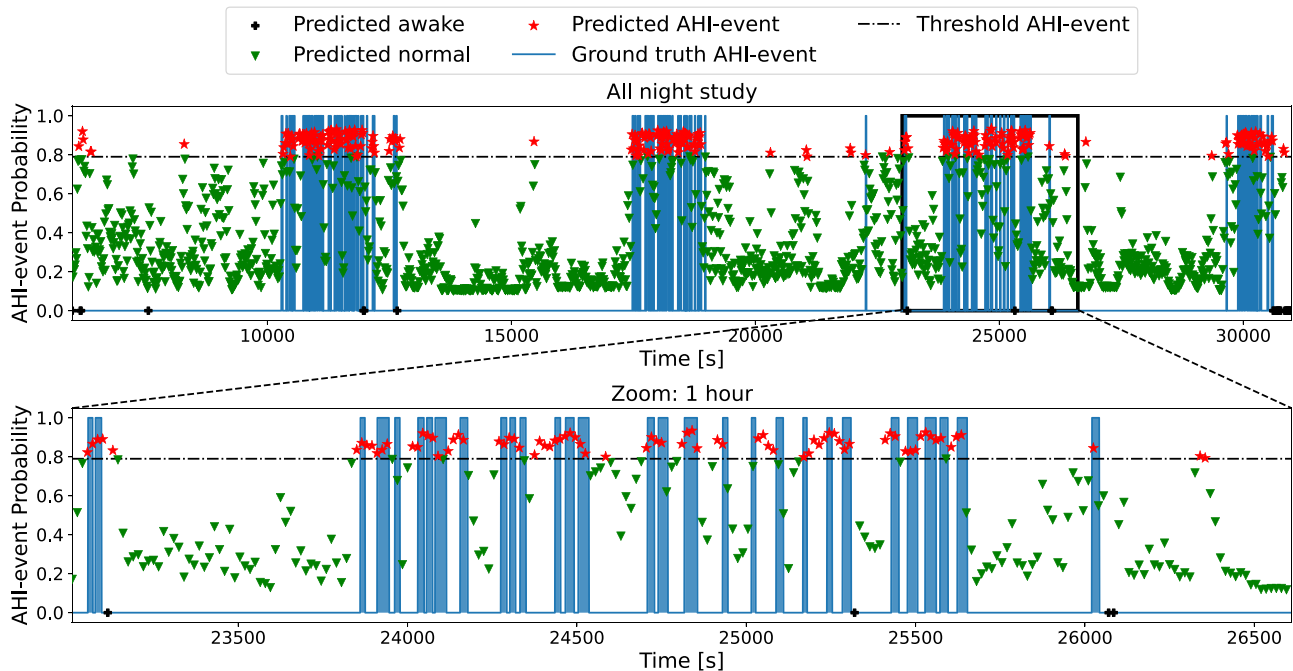
Sensors	F1 AHI
Abdo., Thor., SpO <sub>2</sub>	0.72
Abdo., SpO <sub>2</sub>	0.72
Thor., SpO <sub>2</sub>	0.69
Abdo., Thor.	0.56
SpO <sub>2</sub>	0.63
Abdo.	0.55
Thor.	0.51

single or a combination of easy to use at home sensors: abdominal movement, thoracic movement, and oxygen saturation. We showed that a deep neural network could be effectively trained using an extensive database to extract relevant features from raw physiological data. The detection was enhanced by combining it with a LightGBM classifier, which proved to generalize better than using a CNN on its own, as in ref. 25. Moreover, DRIVEN’s architecture allows for each channel to have its own sampling frequency, and it estimates both AHI events and sleep stages. The best trade-off between model performance and patient comfort was achieved when combining data from only two sensors (abdominal movement and pulse oximetry), yielding an F1-score of 72% at estimating AHI. Hence, DRIVEN can provide an estimation of AHI and detection and segmentation of whole sleep studies to aid physicians make a final diagnosis. For patients, it can guide them to seek medical attention.

Thoracic and abdominal movement sensors each featured very similar performance, since they both measure breathing. Adding SpO<sub>2</sub> significantly improved the performance, which is expected given the role of SpO<sub>2</sub> in the clinical definition of AHI. However, SpO<sub>2</sub> alone struggled to estimate the total sleep time per patient and, for this reason, it should not be used alone. Hence, the combination of SpO<sub>2</sub> and abdominal movement sensors provided an overall good performance both in detecting AHI events and estimating the total sleep time per patient. Supplementary Table 3 reports the performance of DRIVEN on AHI-event classification using all possible combinations of physiological signals available in the considered datasets. RRI data, extracted from ECGs, provided only a marginal improvement combined with other sensors. On its own, it resulted in a poor performance. Nasal airflow was also not included in the main analysis (Figs. 2 and 3) since it is an uncomfortable sensor to sleep with, and provided similar performance as either thoracic movement or abdominal movement sensors.

The reasons for choosing CNNs to automatically search for features over domain knowledge feature extraction methods are mainly due to the variety of channels to evaluate and the large available datasets. The flexibility of deep-learning feature extraction approach allowed us to explore several choices of physiological channels (including all possible combinations). DRIVEN’s architecture allows us to easily and efficiently evaluate new sensors, such as adding a movement sensor if it is available. Several studies have shown that, when the training datasets are large and there is enough computing power, deep learning performs as well or better than other machine-learning algorithms<sup>26</sup>, as demonstrated in different applications ranging from chess games and speech recognition to computer vision<sup>27</sup>. In the particular case of obstructive apnea detection,<sup>28</sup> reported consistently better results using unbiased deep-learning approaches and CNNs for segment classification in PSG studies.

A general limitation, not just of DRIVEN but also any other method without access to EEG data, is detecting hypopneas type 2. This type of hypopnea requires estimating short-term arousal following a drop of at least 30% in nasal airflow, which is particularly challenging to detect from the



**Fig. 4 | Automatic labeling of AHI events for a random patient using two sensors (abdominal movement and SpO<sub>2</sub>).** The blue areas represent true events (zero indicates no event and one indicates an AHI event). The output of DRIVEN is illustrated with symbols that represent, for each 30 s window, the probability of the window to be classified as an AHI event. The windows are colored according to their classification, depending on whether they are above or below the determined

threshold of 0.79. The black crosses represent the segments that were classified as awake, the green triangles the ones classified as not AHI events, and the red stars are the windows classified as AHI events. The second plot zooms in a 1 h segment. Supplementary Fig. 10 increases even further the resolution to a 15 min interval. Supplementary Fig. 11 includes the ground truth labels divided by apnea and different hypopnea types.

three sensors selected in this study (Supplementary Fig. 2). Figure 4 and Supplementary Fig. 10 also show a few false positives and some 30 s segments with a high probability of being AHI events, which in fact are hypopneas type 3 (that are not considered in AHI, Supplementary Fig. 11). Indeed, detecting AHI events (apneas and hypopneas types 1 and 2) is considerably more challenging than simply detecting obstructive sleep apneas (Supplementary Fig. 2), which is the main focus of most available algorithms<sup>28</sup>. A more specific limitation of DRIVEN is its use of 30 s windows, which might treat bursts of multiple short-lived apnea or hypopnea as a single long event. While the smoothing window has merits by effectively filtering noise and capturing longer events, it could lead to underestimating the number of apnea or hypopnea, especially in cases of severe AHI. This effect can be seen in Supplementary Fig. 10, which is a zoomed-in segment of Fig. 4. The smoothing effect caused individual apnea events of bursts to be merged into one event, resulting in a lower AHI estimation. The complex influence of different bursting patterns of apnea events on AHI is also discussed in the literature. For example,<sup>29</sup> demonstrated that very different bursting patterns could lead to the same AHI values. In future studies, we plan to explore new types of AI models, including selective state-space neural networks<sup>30</sup>, to potentially further improve performance and address these limitations. Finally, we investigated classifying sleep or awake from combinations of the three sensors considered in the paper (Supplementary Table 1 and Supplementary Fig. 1). There was a decrease in performance of DRIVEN when using predicted sleep (Fig. 2) instead of the annotated total sleep time per patient obtained from the dataset (Supplementary Fig. 12), as summarized in Supplementary Table 2. In practice, DRIVEN's performance could be further improved by combining these sensors with other easy-to-wear sensors like smartwatches or sensors under the mattress.

This result demonstrates the potential of DRIVEN for an automatic labeling application. Compared to a Polysomnographic Technologist's manual detection and diagnosis of OSA, computer-assisted signal analysis systems can reduce errors related to inter- and intra-operation variability and tiredness induced by the arduous annotating process<sup>31</sup>. Furthermore,

most computer-based analyses can be conducted more cost-effectively and quickly<sup>13</sup>. Overall, DRIVEN could reduce the burden on Polysomnographic Technologists and costs to the healthcare system by making feasible the examination of patients through multiple nights from data that can be collected at home by wearable devices.

## Methods

### Data specification

The polysomnography data is available in European Data Format (EDF)<sup>32</sup>, while the profusion files, in XML format, serve as the annotated label files. These files have been annotated for each sleep stage using the Rechtschaffen & Kales (R&K) criteria<sup>33</sup>. Each profusion file includes annotations for every 30-second sample of sleep stages and instances of sleep problems, with information on initiation and duration time, as well as the sensor used to detect the anomaly. Databases are summarized as follows.

1. The MrOS Sleep Study: MrOS is a substudy of the Men Study of Osteoporotic Fractures. Between 2000 and 2002, a baseline evaluation was performed on 5994 elderly men 65 years of age or older taken from six clinical institutions. Between December 2003 and March 2005, a total of 3135 of these participants were subjected to complete unattended polysomnography and 3 to 5-day actigraphy tests as part of the sleep study. The purpose of the Sleep Study was to determine the extent to which sleep disorders are linked to adverse health outcomes, such as the increased risk of death, fractures, falls, and cardiovascular disease.
2. The MESA Sleep Study: The Multi-Ethnic Study of Atherosclerosis (MESA) is a collaborative 6-center longitudinal investigation of factors associated with the development and progression of subclinical to clinical cardiovascular disease in 6814 black, white, Hispanic, and Chinese American men and women aged 45 to 84 years in 2000–2002. The participants were all between the ages of 45 and 84 at the time of the study. Four follow-up examinations have been performed, one each in the years 2003–2004, 2004–2005, 2005–2007, and 2010–2011. Furthermore, 2237 people participated in a Sleep Exam conducted by

MESA Sleep between 2010 and 2012. This exam included an unattended overnight polysomnogram, wrist-worn actigraphy for seven days, and a sleep questionnaire. The sleep study aims to determine whether there is a correlation between subclinical atherosclerosis and sex, ethnicity, or other demographic differences in sleep and sleep disorders.

3. The SHHS Sleep Study: The Sleep Heart Health Study is a multicenter cohort study conducted by the National Institute of Heart, Lung, and Blood. The purpose of the study was to investigate the effects of sleep-disordered breathing on the cardiovascular system and also on other aspects of a person's health and to determine whether breathing problems that occur during sleep are related to an increased risk of coronary heart disease, stroke, death from all causes, and hypertension. Between November 1995, and January 1998, SHHS Visit 1 research focused on participants who were at least 40 years old and included 6441 men and women. A second polysomnogram, known as SHHS Visit 2, was performed on 3295 participants during the third exam cycle (January 2001 to June 2003).

For the labeling of events, a window is considered AHI event positive if it satisfies one of the AHI conditions (apnea, hypopnea type 1, and hypopnea type 2). To obtain AHI, the total number of AHI events throughout the night is divided by the total number of hours of sleep.

### Data pre-processing and model training

First, data is split by measurement channel (SpO<sub>2</sub>, abdominal movement, thoracic movement, RRI and airflow). For each patient, the initial and final 30 min were removed from the recording as they are considered set up times. Then, for each channel per patient, all signals are standardized and normalized (to speed up training) by calculating the z-scores of 95% of the data points and applying min-max normalization to eliminate the bias of mean and variance of the raw one-dimensional signal<sup>34</sup>. The RRI data was inferred from ECG data by measuring the time difference between heartbeats (from one R peak to the next R peak) using the Pan-Tompkins algorithm<sup>35,36</sup>. Subsequently, each channel's frequency was set to 64Hz, which is the maximum frequency of the measurements, using the cubic spline as the interpolation method<sup>37</sup>; this method is simple to implement and, at the same time, it does not attenuate higher frequency components of the signal. Finally, the signals were segmented into windows of 30 s and labeled as "normal" and "AHI events". Data windows correspond to the same time instance across all channels (Fig. 1a).

Normalized samples of 30 s from each channel are then fed in parallel into distinct deep CNNs. Each CNN is independently trained on a specific channel to classify the AHI events from the normal windows, and then used for automatic feature extraction (Fig. 1b). The training was performed balancing the number of normal and AHI events fed to the network. For this task, the number of normal samples was randomly downsampled to match the amount of positive samples. Once the networks have been trained, features are extracted and concatenated to integrate information from different physiological biomarkers. This is achieved by using the concatenated features to train a LightGBM classifier for binary classification between normal and positive AHI-event windows (Fig. 1c).

The CNNs were trained and cross-validated with 7,753,500 samples, 1,042,652 positives and 6,710,848 negatives, from 5288 PSG recordings in the SHHS database. We use the EfficientNetV2 architecture, a deep CNN with 479 layers developed by Google in 2021<sup>38</sup>. It is a modified and optimized version of EfficientNet<sup>39</sup>, a popular image classification algorithm that won the ImageNet 2019 competition<sup>40</sup>. The architecture used in this paper has been modified to handle unidimensional data and perform binary classification. Each CNN was independently trained using raw unidimensional physiological data from pulse oximetry, RRI, thoracic movement, abdominal movement, or nasal airflow. Categorical cross-entropy was used as the loss function and ADAM as the optimizer<sup>41</sup>. If the validation loss did not decrease after eight consecutive epochs, the training was terminated. Once the networks were trained, the final two layers were removed and the

last global average pooling layer is used to yield 1280 features from each data channel.

### Feature classifier

After the neural networks are trained, and features are extracted, the next step is to concatenate the features extracted by all CNNs and use them for training a LightGBM classifier. In contrast to a large number of other well-known algorithms, such as XGBoost<sup>42</sup> and GBDT<sup>43</sup>, LightGBM employs the classification algorithm for growing trees in a leaf-wise manner rather than in a depth-wise manner. The leaf-wise algorithm can converge significantly faster than the depth-wise growth method, although its growth can be subject to overfitting if the appropriate hyperparameters are not used<sup>23</sup>. We use a random search method within a specified set of parameters to optimize the training and performance of the LightGBM. This allows a fixed number of parameters from a particular distribution to be sampled instead of testing all the values of potential parameters<sup>44</sup>. The output of the feature classifier is the probability that the evaluated segment contains AHI events according to the AHI definition. The training was performed with a balanced dataset. The complete datasets, however, are imbalanced (more negative AHI events than positive ones). Hence, on the validation data, we determined the probability threshold that optimizes the performance metrics of the classifier. This probability threshold was chosen so that the precision and the recall metrics have the same value in the validation dataset (Supplementary Fig. 9a).

### Models and channels comparison

We evaluated longer and shorter window sizes for event classification and AHI estimation on a subset of 1000 randomly selected patients in a 3-fold cross-validation experiment, and selected 30 s windows for their higher performance in the validation set. Additionally, we noticed a delay in the response of SpO<sub>2</sub> signal when a respiratory event occurred; thus, we also evaluated the model performance using different delays in the SpO<sub>2</sub> signal (10, 15, 20, and 25 s delays). For AHI-event classification, we evaluated the feature extraction at different layers of the neural network (model called LGBM-1 for the last layer and LGBM-3 for the third last layer), and the performance of neural networks trained on each channel separately. In summary, we evaluated 9 different channels (abdominal, thoracic, nasal airflow, RRI, SpO<sub>2</sub>, SpO<sub>2</sub>+10s delay, SpO<sub>2</sub>+15s delay, SpO<sub>2</sub>+20s delay, and SpO<sub>2</sub>+25s delay), two different models (LGBM applied to the combination of different sensors extracting the information at the third last layer or at the last layer of each EfficientNetV2), and three different window sizes (10s, 30s, and 80s). Note that every time SpO<sub>2</sub> is used, all its delayed signals are also used as separate channels. The summary of these results are reported in Supplementary Table 3 and the complete set of results is in Supplementary Table 4. These tables show that the 10s windows yield consistently lower performance and RRI signals provide negligible performance improvements. Hence, we now focus on the window sizes of 30 s and 80 s, as well as the channels abdominal, thoracic, and SpO<sub>2</sub>, and use the training data as described below. The complete performance results are reported in Supplementary Table 5, which also includes the estimation of AHI. On the validation dataset, these tables show that a larger window size leads to better classification of AHI events, whereas a lower window size results in better AHI estimation. Larger window sizes are less suited for counting the AHI events needed for AHI estimation as multiple AHI events are counted as a single detection. Since most AHI events last 30 s or less, the 30 s window presents a good trade-off between event classification and AHI estimation. Hence, we finally selected the 30 s windows because the objective of this paper is on AHI estimation. The 30 s windows also allow us to do a better whole sleep study segmentation.

### End-to-end CNN

DRIVEN's architecture consists of a combination of CNNs and LightGBM. Here, we compare its performance with that of a single CNN that inputs all channels and then performs the final classification. For this model, we used the same NN, EfficientNetV2, and changed the input layer to accept a vector

for each channel. As in the previous section, we compared the AHI-event classification as well as the AHI classification of patients with the same 1000 randomly selected patients. Supplementary Table 4 shows that the validation results are quite similar, with the end-to-end CNNs having a slightly lower performance compared with the combination of CNNs and LightGBM. However, the performance of the end-to-end CNNs drop significantly in the test datasets, which are also detailed in this spreadsheet. The end-to-end CNNs tend to overfit the training data and the combination of CNNs and LightGBM shows better generalization across datasets. This result has also been observed by others (see, for example,<sup>25</sup>).

### Dataset partition into train, validation, and test

We performed three different trials with the complete dataset: (1) SHHS and MROS patients as training and validation sets in an 80–20% random split and MESA patients as test set, (2) SHHS patients as training and validation sets in an 80–20% random split and MESA and MROS as test set, and (3) SHHS, MROS and MESA patients as training, validation and test sets in a 60–20–20% random split. Because there are patients in SHHS and MROS databases that have two recordings, to avoid data leakage we were careful to always place them in the same partition. The performance metrics for validation and test sets of each trial can be found in Supplementary Table 5. All the results and figures presented in the paper are from trial 2, as we only trained and validated with patients from the SHHS database and tested the model's performance to close the domain gap when applying our algorithm to the other two independent databases.

### Sleep/awake classification

The calculation of AHI requires the total time the patient has slept, excluding awake times throughout the night. Hence, it is important to discriminate between sleep and awake intervals. The profusion files per patient were already labeled with this information, which we could use to compute the total sleeping time. This information can also be estimated from other sensors like photoplethysmograms (PPG)<sup>45</sup> or accelerometers. Here, to provide a standalone solution, we investigated classifying sleep or awake from combinations of only the available three sensors considered in the paper. We used the same architecture (CNN+LGBM) and model comparison methodologies to estimate if a 30 s window is sleep or awake. Additionally, we investigated adding an extra sensor, RRI, already discarded for AHI-event prediction. The results are presented in Supplementary Table 1 and Supplementary Fig. 1. The performance is high for most combinations of sensors. The exception is SpO<sub>2</sub> alone, which struggles to correctly differentiate between sleep and awake windows. Adding an RRI sensor does not improve the sleep prediction performance either with the best classifiers or alone. This sleep prediction, with a threshold of 0.5 to classify sleep or awake, is then used for the AHI estimation of each patient.

### AHI and severity class estimation from positive events

As the data were imbalanced, we determined the final classification threshold as the point of intersection of the precision and the recall curves on the validation data (Supplementary Fig. 9a). Then, for each patient in the validation dataset, we went through the overnight study classifying all 30 s windows sampled every 15 s (hence, consecutive windows overlap by 15 s). Each 30s window was first classified as awake or sleep, and the sleep segments were further classified into positive AHI event or negative by using the previously defined threshold (Fig. 4). Next, we calculated the ratio between the number of positive AHI events over the total sleep segments. With these ratios, one per patient, we fitted a linear regression to the value of AHI of these patients. In this regression, we used weights to account for the imbalance of the number of patients in different AHI classes. For the test set, we used the threshold selected before (Supplementary Fig. 9b), and obtained the AHI event vs sleep windows ratio as before. Finally, we used the previously determined linear regression factors to obtain the AHI of the test patients. Recall that the patient severity class was determined following the rules provided by the American Academy of Sleep Medicine<sup>6</sup>.

### Performance metrics

The performance of the models was evaluated based on the following metrics. Let true positive (TP) represent the number of AHI events that are accurately predicted. True negative (TN) represents the number of normal events that are accurately predicted. False negative (FN) represents the number of AHI events incorrectly predicted as normal events, and false positive (FP) represents the number of normal events incorrectly predicted as AHI events. The accuracy of the model is given by  $(TP + TN)/(TP + TN + FP + FN)$ , indicating the probability of correctly identifying normal and AHI events; the recall is given by  $TP/(TP + FN)$ , indicating the probability of identifying AHI events; the specificity is given by  $TN/(TN + FP)$ , indicating the probability of detecting normal events; the precision is given by  $TP/(TP + FP)$ , indicating the ratio between detected AHI events and all predicted AHI-event cases; and the F1-score  $2(\text{precision} \times \text{recall})/(\text{precision} + \text{recall})$ , indicating the harmonic mean of precision and recall. For every method, the receiver operating characteristic curve, and the precision-recall curve were obtained for each possible classification threshold. The area under both of these curves allows us to compare the global performance of each architecture. Finally, the patient severity classification results were evaluated by assessing the F1 score of each class and then averaging them.

### Data availability

The data was provided by the National Sleep Research Resource and are publicly available on request at <https://sleepdata.org/>.

### Code availability

Data pre-processing and segmentation were implemented using MATLAB software. The neural network was implemented on the Keras framework with Tensorflow backend on Python 3.7. The codes are available in GitHub at <https://github.com/LCSB-SCG/DRIVEN/>.

Received: 6 February 2023; Accepted: 22 May 2024;

Published online: 01 June 2024

### References

1. Senaratna, C. V. et al. Prevalence of obstructive sleep apnea in the general population: a systematic review. *Sleep. Med. Rev.* **34**, 70–81 (2017).
2. Riha, R. L. Defining obstructive sleep apnoea syndrome: a failure of semantic rules. *Breathe* **17**, 210082 (2021).
3. Mannarino, M. R., Di Filippo, F. & Pirro, M. Obstructive sleep apnea syndrome. *Eur. J. Intern. Med.* **23**, 586–593 (2012).
4. Punjabi, N. M. The epidemiology of adult obstructive sleep apnea. *Proc. Am. Thorac. Soc.* **5**, 136–143 (2008).
5. McNicholas, W. T. Diagnosis of obstructive sleep apnea in adults. *Proc. Am. Thorac. Soc.* **5**, 154–160 (2008).
6. Berry, R. B. et al. Rules for scoring respiratory events in sleep: update of the 2007 AASM manual for the scoring of sleep and associated events: deliberations of the sleep apnea definitions task force of the American Academy of Sleep Medicine. *J. Clin. Sleep. Med.* **8**, 597–619 (2012).
7. Hutchison, K. N., Song, Y., Wang, L. & Malow, B. A. Analysis of sleep parameters in patients with obstructive sleep apnea studied in a hospital vs. a hotel-based sleep center. *J. Clin. Sleep. Med.* **4**, 119–122 (2008).
8. Levy, J., Álvarez, D., Del Campo, F. & Behar, J. Deep learning for obstructive sleep apnea diagnosis based on single channel oximetry. *Nat. Commun.* **14**, 4881 (2023).
9. Ramachandran, A. & Karuppiah, A. A survey on recent advances in machine learning based sleep apnea detection systems. *Healthcare* **9**, 914 (2021).
10. Mendonca, F., Mostafa, S. S., Ravelo-Garcia, A. G., Morgado-Dias, F. & Penzel, T. A review of obstructive sleep apnea detection approaches. *IEEE J. Biomed. Health Inform.* **23**, 825–837 (2018).

11. Mostafa, S. S., Mendonça, F., G. Ravelo-García, A. & Morgado-Dias, F. A systematic review of detecting sleep apnea using deep learning. *Sensors* **19**, 4934 (2019).
12. Yumino, D. et al. Differing effects of obstructive and central sleep apneas on stroke volume in patients with heart failure. *Am. J. Respiratory Crit. Care Med.* **187**, 433–438 (2013).
13. Faust, O., Hagiwara, Y., Hong, T. J., Lih, O. S. & Acharya, U. R. Deep learning for healthcare applications based on physiological signals: a review. *Computer Methods Prog. Biomedicine* **161**, 1–13 (2018).
14. Tang, J., Alelyani, S. & Liu, H. Feature selection for classification: a review. *Data Classification: Algorithms Appl* **37** (2014).
15. Ganapathy, N., Swaminathan, R. & Deserno, T. M. Deep learning on 1-d biosignals: a taxonomy-based survey. *Yearb. Med. Inform.* **27**, 098–109 (2018).
16. Ribeiro, A. et al. Automatic diagnosis of the 12-lead ECG using a deep neural network. *Nat. Commun.* **11**, 1760 (2020).
17. Gavidia, M. et al. Early warning of atrial fibrillation using deep learning. *Patterns* **5**, 100970 (2024).
18. Chen, X. et al. Racial/ethnic differences in sleep disturbances: the multi-ethnic study of atherosclerosis (mesa). *Sleep* **38**, 877–888 (2015).
19. Blackwell, T. et al. Associations between sleep architecture and sleep-disordered breathing and cognition in older community-dwelling men: the osteoporotic fractures in men sleep study. *J. Am. Geriatrics Soc.* **59**, 2217–2225 (2011).
20. Quan, S. F. et al. The sleep heart health study: design, rationale, and methods. *Sleep* **20**, 1077–1085 (1997).
21. Zhang, G.-Q. et al. The national sleep research resource: towards a sleep data commons. *J. Am. Med. Inform. Assoc.* **25**, 1351–1358 (2018).
22. Arnold, J., Boucher, J., Mobley, D., Nawabit, R. & Redline, S. *SRC Manual of Operations and Scoring Rules* (MESA PSG Sleep Reading Center, 2014).
23. Ke, G. et al. Lightgbm: A highly efficient gradient boosting decision tree. *Adv. Neural Inform. Process. Syst.* **30** (2017).
24. Mohammed, R., Rawashdeh, J. & Abdullah, M. Machine learning with oversampling and undersampling techniques: overview study and experimental results. In *2020 11th Int. Conference On Information And Communication Systems (ICICS)* 243–248 (IEEE, 2020).
25. Bentéjac, C., Csörgő, A. & Martínez-Muñoz, G. A comparative analysis of gradient boosting algorithms. *Artif. Intell. Rev.* **54**, 1937–1967 (2021).
26. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **512**, 436–444 (2015).
27. Sutton, R. The bitter lesson. *Incomplete Ideas (blog)* **13**, 38 (2019).
28. Moridian, P. et al. Automatic diagnosis of sleep apnea from biomedical signals using artificial intelligence techniques: Methods, challenges, and future works. *WIREs Data Min. Knowl. Discov.* **12**, e1478 (2022).
29. Chen, S., Redline, S., Eden, U. & Prerau, M. Dynamic models of obstructive sleep apnea provide robust prediction of respiratory event timing and a statistical framework for phenotype exploration. *Sleep* **45**, zsac189 (2022).
30. Gu, A. & Dao, T. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. Preprint at <https://arxiv.org/abs/2312.00752> (2023).
31. Alvarez-Estevez, D. & Moret-Bonillo, V. Computer-assisted diagnosis of the sleep apnea-hypopnea syndrome: a review. *Sleep Disorders* **2015** (2015).
32. Kemp, B. & Olivan, J. European data format ‘plus’(edf+), an edf alike standard format for the exchange of physiological data. *Clin. Neurophysiol.* **114**, 1755–1761 (2003).
33. Rechtschaffen, A & Kales, A. A manual for standardized terminology, techniques and scoring system for sleep stages in human subjects. *Public Health Service, US Government Printing Office, Washington DC.* (1968).
34. Huang, L. et al. Normalization techniques in training DNNs: Methodology, analysis and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence.* **45**, 10173–10196 (2023).
35. Pan, J. & Tompkins, W. J. A real-time QRS detection algorithm. *IEEE Trans. Biomed. Eng.* **BME-32**, 230–236 (1985).
36. Sedghamiz, H. Matlab implementation of pan tompkins ECG QRS detector. Code Available at File Exchange Site of Mathworks (2014).
37. Keys, R. Cubic convolution interpolation for digital image processing. *IEEE Trans. Acoust., Speech, Signal Process.* **29**, 1153–1160 (1981).
38. Tan, M. & Le, Q. EfficientNetV2: Smaller models and faster training. In *Int. Conference on Machine Learning. PMLR* 10096–10106 (2021).
39. Tan, M. & Le, Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In *Int. Conference on Machine Learning. PMLR*, 6105–6114 (2019).
40. He, K., Girshick, R. & Dollár, P. Rethinking imagenet pre-training. In *Proc. IEEE/CVF International Conference on Computer Vision. ICCV*, 4918–4927 (2019).
41. Jais, I., Ismail, A. & Nisa, S. Adam optimization algorithm for wide and deep neural network. *Knowl. Eng. Data Sci.* **2**, 41–46 (2019).
42. Chen, T. et al. Xgboost: extreme gradient boosting. *R. Package Version 0. 4-2* **1**, 1–4 (2015).
43. Ye, J., Chow, J.-H., Chen, J. & Zheng, Z. Stochastic gradient boosted distributed decision trees. In *Proc. 18th Acm Conference on Information And Knowledge Management*, 2061–2064 (2009).
44. Bergstra, J. & Bengio, Y. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* **13** (2012).
45. Kotzen, K. et al. SleepPPG-Net: A Deep Learning Algorithm for Robust Sleep Staging From Continuous Photoplethysmography. *IEEE J. Biomed. Health Inform.* **27**, 924–932 (2023).
46. Varrette, S., Bouvry, P., Cartiaux, H. & Georgatos, F. Management of an Academic HPC Cluster: The UL Experience. In *Proc. 2014 Int. Conference on High Performance Computing & Simulation. IEEE*, 959–967 (2014).

## Acknowledgements

The authors acknowledge support from the Luxembourg National Research Fund (FNR) through grants PRIDE15/10907093/CriTICS, AFR/17022833 and INTER/DFG/21/15020234, the latter is co-funded by the Deutsche Forschungsgemeinschaft (DFG) project number 458610525. High performance computing experiments for model training and evaluation presented in this report were carried out using the HPC facilities of the University of Luxembourg<sup>6</sup>.

## Author contributions

J.G. conceptualized the research; M.E.G., J.G., G.R., and A.N.M. developed the AI model; G.R., M.E.G., and B.B. implemented the codes and validated the results; G.R., M.E.G., J.G., B.B., A.H., and A.N.M. analyzed the results; J.G. and A.H. supervised the work; all authors wrote and revised the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41746-024-01139-z>.

**Correspondence** and requests for materials should be addressed to Jorge Goncalves.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024