

Ubi-SleepNet: Advanced Multimodal Fusion Techniques for Three-stage Sleep Classification Using Ubiquitous Sensing

BING ZHAI, School of Computing, Newcastle University, UK

YU GUAN, School of Computing, Newcastle University, UK

MICHAEL CATT, Population Health Sciences Institute, Newcastle University, UK

THOMAS PLÖTZ, School of Interactive Computing, Georgia Institute of Technology, USA

Sleep is a fundamental physiological process that is essential for sustaining a healthy body and mind. The gold standard for clinical sleep monitoring is polysomnography (PSG), based on which sleep can be categorized into five stages, including wake/rapid eye movement sleep (REM sleep)/Non-REM sleep 1 (N1)/Non-REM sleep 2 (N2)/Non-REM sleep 3 (N3). However, PSG is expensive, burdensome and not suitable for daily use. For long-term sleep monitoring, ubiquitous sensing may be a solution. Most recently, cardiac and movement sensing has become popular in classifying three-stage sleep, since both modalities can be easily acquired from research-grade or consumer-grade devices (e.g., Apple Watch). However, how best to fuse the data for greatest accuracy remains an open question. In this work, we comprehensively studied deep learning (DL)-based advanced fusion techniques consisting of three fusion strategies alongside three fusion methods for three-stage sleep classification based on two publicly available datasets. Experimental results demonstrate important evidences that three-stage sleep can be reliably classified by fusing cardiac/movement sensing modalities, which may potentially become a practical tool to conduct large-scale sleep stage assessment studies or long-term self-tracking on sleep. To accelerate the progression of sleep research in the ubiquitous/wearable computing community, we made this project open source, and the code can be found at: <https://github.com/bzhai/Ubi-SleepNet>.

CCS Concepts: • **Computing methodologies** → **Neural networks**; • **Human-centered computing** → **Ubiquitous and mobile computing**.

Additional Key Words and Phrases: Three Sleep Stages, Sleep Monitoring; Deep Learning; MESA; Apple Watch, Wearable, Heart Rate, Heart Rate Variability, Multimodal Fusion, Ubiquitous Sensing, Neural Networks

ACM Reference Format:

Bing Zhai, Yu Guan, Michael Catt, and Thomas Plötz. 2021. Ubi-SleepNet: Advanced Multimodal Fusion Techniques for Three-stage Sleep Classification Using Ubiquitous Sensing. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 33 pages. <https://doi.org/xxxxx/1xxx.xxxxxx6>

1 INTRODUCTION

Human beings spend about one-third or more of their lives sleeping. It is a physiological process essential to maintain life and health [71], and it plays a major role in repairing the body tissues, removing toxic metabolic waste products from the brain, consolidating memory, restoring energy and enhancing immune defence [7, 10, 23, 57, 65]. Long-term lack of sleep and/or poor-quality sleep is likely to increase the risk of obesity, diabetes, and heart and blood vessel (cardiovascular) disease [38, 57]. Accurate and long-term sleep monitoring using ubiquitous sensing

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/xxxxx/1xxx.xxxxxx6>

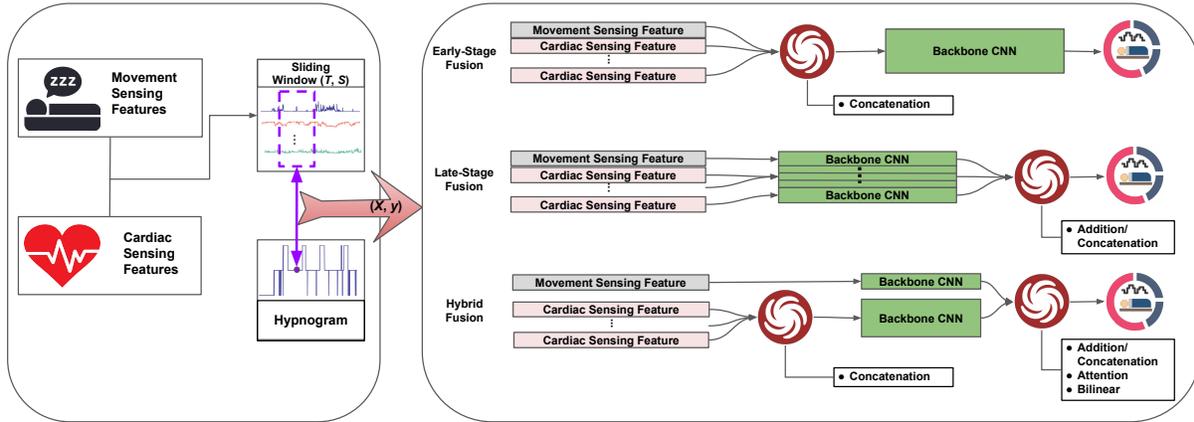


Fig. 1. An overview of the three-stage sleep classification system. Features were extracted for each sleep epoch (30s). The sliding window method divides the sleep data into multiple segments with window length T and stride S . In our study, we set $T = 101$, and $S = 1$. The hypnogram represents the stages of sleep over time in each sleep epoch. Three fusion strategies and three fusion methods were studied. The prediction process was performed for each sleep epoch.

technology is increasingly vital to the understanding of human health and is becoming an active area of health research.

The gold standard for clinical sleep monitoring is polysomnography (PSG), which requires subjects to wear multi-channel sensors on the body [35]. PSG recording can be classified into five stages, i.e., wake, rapid eye movement sleep (REM sleep) and three-types of non-rapid eye movement sleep (NREM sleep) including N1, N2 and N3. According to the American Academy of Sleep Medicine (AASM) rules [33], each stage lasts 30 seconds (i.e., a sleep epoch). However, it is expensive and burdensome, which is impractical for long-term sleep monitoring.

For such monitoring, many ubiquitous sensing approaches were studied, including actigraphy [4], smart watches [74], WiFi [82], bed sensors [55] and radio signals [32], etc. Among them, in terms of reliability and usability, cardiac and movement (upper limb) sensing are considered promising modalities. They can be easily collected from lightweight research/consumer-grade devices (e.g., Apple Watch [74]). Based on cardiac and movement sensing, previous work studied a number of machine learning and DL approaches on sleep monitoring tasks from two-stage (wake/sleep) to five-stage (wake/REM/N1/N2/N3) sleep classification [84]. Their initial results suggested the feasibility of using these two modalities for three-stage (wake/REM/NREM) sleep classification, and their findings corroborate sleep physiology studies [16, 50], where NREM sleep was not deemed to be easily separated into N1/N2/N3 without employing EEG signals.

The easy-to-collect nature of cardiac and movement sensing provided a salable method for large-scale and long-term sleep monitoring. Longitudinal sleep monitoring with accurate details in (three) sleep stages is meaningful to health and medical research. Deep NREM sleep (or slow wave sleep - SWS) is known to be the most “restorative” sleep stage, which controls hormonal changes that affect glucose regulation [19]. Long-term reduction in NREM sleep may adversely affect glucose homeostasis and increase the risk of type 2 diabetes [64]. The REM sleep dysregulation has played a central role in depression and Parkinson’s studies [3, 39]. For instance, reduced REM sleep latency, along with increased REM sleep duration and REM sleep density, have been considered to be an objective indicator of depressive disorder and inversely correlated to its severity [40, 68, 75]. The increased health research density in digital phenotypes by using inexpensive, mass-produced consumer wearables demand reliable

algorithms that can classify sleep stages in longitudinal settings [70]. Beyond health and clinical monitoring, the sleep monitoring has also been welcomed by self-trackers in the past decades [63].

To the best of our knowledge, only two previous studies [74, 84] have adopted cardiac and movement sensing for three-stage sleep classification based on publicly accessible sleep datasets. Both works used very basic multimodal fusion techniques (i.e., feature concatenation [74, 84]) in neural networks, which tested the feasibility of classifying three-stage sleep. However, the models used in these benchmark studies did not achieve good performance, due to overestimating NREM sleep time and underestimating wake time. The simple fusion technique may not fully utilize the advantages of multimodal data, especially in heterogeneous multimodal data scenarios. For instance, movement sensing can achieve better classification performance in sleep/wake tasks compared with the use of cardiac sensing alone [84]. But movement sensing alone is incapable of discerning three-stage sleep. Given that, it is desirable to explore the advanced fusion techniques to further boost the performance.

In this work, we firstly systematically studied three fusion strategies for three-stage sleep classification, including early-stage fusion, late-stage fusion and hybrid fusion, to answer the question, *"At what stage should the cardiac and movement sensing representation be merged?"*

Secondly, we employed three fusion methods (simple operations, attention mechanism, tensor methods) to answer the question, *"How to better combine cardiac and movement sensing representations?"* The simple baseline operations (concatenation and addition) as well as the advanced fusion methods (the attention mechanism-based method [79] and bi-linear pooling-based method [45]) were studied. The pipeline of our system is demonstrated in Figure 1.

These fusion techniques were comprehensively evaluated on two public datasets which are the Apple Watch dataset [74] and the Multi-Ethnic Study of Atherosclerosis (MESA) dataset [14, 84, 86]. The Apple Watch dataset includes cardiac and movement signals collected from consumer-grade devices from a cohort of 31 young and healthy adults. For the MESA dataset, only the cardiac and movement sensing signals were used, which can be acquired from research-grade devices. The dataset consists of 1743 subjects from the aging population.

For these two representative datasets, our results suggested that three-stage sleep classification can be reliably achieved by employing advanced fusion techniques on the cardiac and movement sensing data, which can be easily acquired from consumer/research-grade devices. Several models developed in our study achieved the state-of-the-art performance for three-stage sleep classification. We also evaluated the module parameter size and its corresponding inference time, which may play a vital role in ubiquitous computing applications.

Moreover, we also investigated a visualization method to explore the decision-making process of the multimodal fusion model for three-stage sleep classification. The exploratory user research demonstrated that the gradient class activation map (Grad-CAM) [66] based sleep data visualization can be understood and used by humans, which facilitates the transparency of using DL in sleep health research.

This work contributes to the long-term non-intrusive three-stage sleep monitoring solution that may be deployed with mass-produced and inexpensive consumer-grade wearables, which may potentially be used for large-scale population-based sleep health studies and long-term sleep self-tracking.

2 SLEEP MONITORING AND UBIQUITOUS SENSING

2.1 Clinical Sleep Monitoring and Sleep Physiology

2.1.1 Polysomnography and Sleep Stages. Traditionally, gold-standard human sleep assessment was conducted in laboratory settings using polysomnography (PSG) which commonly involved electroencephalography (EEG), electromyography (EMG) and electrooculography (EOG). Together they facilitated the measurement of brain activity, alongside both muscle and eye movement [8]. PSG usually requires a sleep laboratory and a sleep technician in a controlled environment. Multiple skin electrodes will be placed on the subject's body and the

set-up procedure will take 45-60 minutes [58]. It is impractical to measure sleep using this method for more than two consecutive nights as participants feel burdened.

Sleep stages consist of wakefulness (Wake), REM sleep and NREM sleep, three vigilance states for humans [33]. NREM sleep can be further subdivided into three stages: light sleep stage 1 (N1), light sleep stage 2 (N2) and deep sleep (N3) [9]. The alternate appearance of NREM sleep and REM sleep constitutes a sleep cycle. A healthy person typically has 4-5 sleep cycles in one night [8]. As the sleep cycle increases, the proportion of REM sleep increases, and the proportion of NREM sleep decreases. Five sleep stages can be distinguished by PSG, i.e., by analyzing the characteristics of EEG, EMG and EOG.

2.1.2 Actigraphy and Sleep-Wake Monitoring. Actigraphy is a valid method for detecting sleep/wake and is commonly used for ambulatory monitoring of sleep time or rhythms [4, 35]. The actigraphy equipment is a type of wearable wristband that consists of various sensors that can monitor the light-off time and the movement of the limbs (using an accelerometer). However, it is limited to monitoring sleep-wake as the actigraphy data may not contain sufficient information to discern sleep stages.

2.1.3 Cardiac Activities and Sleep Physiology. Cardiovascular autonomous control plays an essential role in sleep, and it will be different when transitioning to different sleep stages. The modulation of the autonomic nervous system (ANS) regulates cardiovascular functions during sleep onset and sleep stages [46, 50]. Heart rate variability (HRV) analysis is a classical tool for the ANS analysis. Research on HRV in sleep stages noted that REM sleep was characterized by a likely sympathetic predominance, while NREM sleep followed an opposite trend [13, 47, 51]. The transition between wake, NREM and REM sleep is accompanied by changes of several HRV characters, such as the HR, Low-Frequency (LF) power, High-Frequency (HF) power and LF/HF ratio [13, 16, 47].

Not all sleep stages are associated with brain activity. A study conducted by Desseilles et al. [16] through HRV and brain imaging analysis found close connectivity between autonomic cardiac modulations and the activity of certain brain areas during REM sleep. There is no conclusive connectivity between the brain and cardiac activity during NREM sleep. Therefore, it may be not be easy to discern each NREM sleep stage accurately, without EEG signals.

2.2 Ubiquitous Sensing Techniques for Sleep Monitoring

2.2.1 Wireless Sleep Monitoring. In recent years, the wireless technologies showed the potential of sleep monitoring, e.g., ballistocardiograph (BCG, with an example consumer product: Beddit™), which can monitor heart rate wirelessly [24]. However, wireless-based approaches face some challenges when deploying them in clinical research owing to, a) the non-standardised measurement methods; b) the lack of precise understanding of the physiological origins that influences the signal waveform; c) comparatively low reliability and specificity of these signals to the existing clinical methods (for example, WiFi signals may be scattered by multiple subjects), which may hamper its wide applications in health and medical research [24].

2.2.2 Mobile and Miniaturized EEG. The development of mobile EEG systems (e.g., Emotiv™) reduced the time spent on the electrodes installation process and improved the motion tolerance during the recording procedure [17]. The main disadvantage of these devices is that the electrodes embedded in the cap are visible, and the form factor limitations prevent comfortable, continuous, and long-term sleep monitoring. Furthermore, once the electrode gel dries, the signal quality decreases, and the gel leaves residues [11], which may impact the wearing experience. The lightweight headband EEG (e.g. Sleep Profiler™ and Dreem™) can monitor sleep in a natural environment. But it takes extra effort each time of wearing to adjust the equipment position to reduce the skin impedance to an acceptable level [44]. In addition, compared to smartwatches, these devices are usually expensive. The recent development of ear EEG improved the wearing comfort and reduced electrode set-up time. Several

studies demonstrated reasonable results of sleep stage classification using these prototypes [48, 49]. However, these ear EEG devices commonly adopt around-ear or/and in-ear style and are made of silicone materials, which offer a bearable wearing experience, making them less popular than the mass-produced wearables [48].

2.2.3 Consumer and Research-grade Wearables for Sleep Monitoring. Many leading consumer products such as Fitbit™ and Xiaomi™ band provide sleep stages tracking services. However, these consumer products commonly lack minimal validation, with poor algorithm transparency on data processing/sleep stage classification, resulting in these devices being precluded in clinical, research, or occupational settings [36]. Nevertheless, another consumer product, Apple Watch, provides access to the accelerometer data and heart rate data, making it feasible to develop an algorithm for sleep health studies and self-trackers.

Consumer-grade wearable devices with diverse modalities offer a potential solution to ambulatory sleep tracking. Such sensors provide valuable, inexpensive, unobtrusive measurement tools to collect biological signals. Many of these wearables can communicate with smartphones, facilitating data collection and storage during large-scale population studies. Therefore, exploring the use of consumer-grade wearables in sleep and health studies becomes prevalent as the HR/HRV data and movement sensing data are generally available on these wearables [74].

The sleep stage classification based on ECG/PPG signals has also been investigated by [20, 21]. The results demonstrated promising performance. However, PPG data is generally unavailable on many consumer wearables, and ECG requires the skin electrodes to be placed near the heart. Collecting these raw signals may require research-grade wearables (e.g., Empatica™ E4), which demand additional financial costs for daily sleep monitoring.

Several previously published studies demonstrated that using HR/HRV features and movement sensing together could discern three sleep stages and achieved promising results [28, 74, 84]. Heterogeneous modalities may carry supplementary information for sleep stage classification. There is still much to be understood regarding how to construct this fusion architecture and which fusion method will be the most effective for sleep-stage classification. Our work adds to this knowledge. Exploring multimodal fusion strategies and methods to better integrate different physiological signals is of great significance for health research and self-monitoring of sleep using ubiquitous computing technology.

3 ADVANCED FUSION TECHNIQUES FOR THREE-STAGE SLEEP CLASSIFICATION

In this section, we will first discuss the current progress of multimodal fusion strategies and methods and their applications in sleep monitoring. We then present our study structure, followed by a technical description of three fusion strategies (early-stage, late-stage and hybrid fusion) and three methods (simple operation, attention mechanism and tensor-based method)

3.1 Overview of Multimodal Fusion

Multimodal fusion in machine learning has been extensively studied in pattern recognition applications, such as in image and video captioning[80], visual question answering [79], audio-visual speech recognition [2] and affect recognition [34]. In the field of ubiquitous computing, multimodal fusion has also been adopted for human activity recognition [43], sleep stage classification [84], fatigue assessment [6] and person identification [15]. The simple concatenation method was commonly adopted in these studies to combine the raw inputs or combine the representations obtained from the pre-trained model of each modality [61]. Other researchers explored more advanced fusion methods, such as the attention-based fusion scheme for human activity recognition [43]

For monitoring sleep, several previous works achieved promising results for sleep stage classification by concatenating multimodal intermediate features and feeding them into DL models [28, 84]. However, these studies focus on the choice of modalities rather than the fusion techniques. Different modalities may contain complementary information. It is difficult to explicitly identify the best suitable cross-modal fusion architectures.

In terms of the movement sensing and cardiac sensing, they are different in signal-to-noise ratio, data generation process and measurement frequency. Moreover, the activity count is better in sleep/wakefulness classification, but it is difficult to discern different sleep stages [4]. For healthy adults, the difference in heart rate variability between REM sleep and wake is less than the difference in NREM and REM sleep [16].

The choice of fusion strategy and fusion method may thus influence the model classification performance. In recent years, the DL-based computational models have outperformed shallow machine learning models for sleep stage classifications, not only on unimodal data, but also on the multimodal data [54, 59, 84]. Therefore, this work will only focus on the multimodal fusion techniques based on DL networks.

3.2 Problem Statement

Based on the movement sensing and cardiac sensing data, the goal of this work is to comprehensively study how to use the advanced fusion techniques to reliably classify three-stage sleep. As demonstrated in Fig. 1, we adopt a sliding window method with window size T and stride S to segment sleep recordings into frames. In each frame, we could either extract the handcraft features (e.g. heart rate features that were deemed to be intermediate / mid-level features) from each sleep epoch that can provide physiologically meaningful features to the model [20, 84], or we could use neural networks to extract the deep features. The time step t represents one sleep epoch (i.e. every 30 seconds). Given that, we aim to map the data in a sliding window to a sleep stage that corresponds to the center point of the window (e.g., the purple point in the hypnogram in Fig. 1).

Suppose the i th frame-wise time-series input data for cardiac sensing can be denoted as $\mathbf{X}_{car}^{(i)} \in \mathbb{R}^{C_{car} \times T}$, where C_{car} denotes the number of features/input channels and T denotes the sliding window length. For movement sensing, the input data can be denoted as $\mathbf{X}_{mov}^{(i)} \in \mathbb{R}^{C_{mov} \times T}$. The details of feature extraction will be introduced in Section 4. The goal of deep multimodal fusion is to determine a multilayer neural network $f(\cdot)$ whose output $\hat{y}^{(i)}$ is expected to be the same as the target $y^{(i)}$ as much as possible for each sample $(\mathbf{X}_{mov}^{(i)}, \mathbf{X}_{car}^{(i)})$. This can be implemented by minimizing the empirical loss \mathcal{L} for classification denoted as:

$$\min_f \frac{1}{N} \sum_{i=1}^N \mathcal{L} \left(\hat{y}^{(i)} = f(\mathbf{X}_{mov}^{(i)}, \mathbf{X}_{car}^{(i)}), y^{(i)} \right) \quad (1)$$

3.3 Fusion Strategy

Traditional fusion strategies include feature level fusion (e.g., [62, 76, 77]), score-level fusion (e.g., [26]) or decision-level fusion (e.g., [25]). In the end-to-end DL era, the boundary between multimodal representation and fusion has been blurred. Representation learning is interlaced with classification (or regression) objectives. Nevertheless, the fusion strategy for DL models may still be carried out in three stages, such as early fusion, late fusion and hybrid fusion [85].

Fusion at different stages may influence the results of representation learning. For example, the early and late fusion may inhibit intra-modal or inter-modal interaction [85]. Neverova et al. noted that highly correlated modalities should be fused together [52]. Hazirbas et al. demonstrated that the performance of fusion is highly affected by the choice of which layer to fuse [29]. For sleep stage classification, the way to fuse heterogeneous intermediate features is worthy of exploration. In this study, we want to gain a comprehensive understanding at what stage we should fuse these inputs to achieve the most performance improvements on the three-stage sleep classification task. We adopted three commonly used fusion strategies, including early-stage fusion, late-stage fusion and hybrid fusion, as shown in Figure 1.

3.3.1 Early-stage Fusion. In the early-stage fusion, data from different modalities (e.g., intermediate features) are concatenated (stacked) in the input stage. It is popular because of its simplicity, yet it is sub-optimal [60]. Early-stage fusion firstly concatenates the cardiac (denoted as subscript *car*) and movement (denoted as subscript

act) sensing data then feeds them into neural networks h to make a corresponding prediction.

$$\hat{y}^{(i)} = h(\text{Concatenate}(\mathbf{X}_{car}^{(i)}, \mathbf{X}_{mov}^{(i)})). \quad (2)$$

where Concatenate is the matrix concatenation function.

3.3.2 Late-stage Fusion. Late-stage fusion is another prevalent way to fuse (high-level) representation from multiple sources. This fusion strategy allows high-level representations to have better intra-modal coherence. Late-stage fusion processes each modality's c th channel input data with a network q and then combines all their high-level representations via an aggregation operation followed by the classification layers. It is denoted as:

$$\hat{y}^{(i)} = \varphi(\text{Agg}(q(\mathbf{x}_{mov,1}^{(i)}), \dots, q(\mathbf{x}_{mov,C_{mov}}^{(i)}), q(\mathbf{x}_{car,1}^{(i)}), \dots, q(\mathbf{x}_{car,C_{car}}^{(i)}))) \quad (3)$$

where Agg is the aggregation function and φ denotes the classifier (e.g. fully connected layers), and $\mathbf{x}^{(i)} \in \mathbb{R}^{1 \times T}$ and T is the window length. The cardiac intermediate features are denoted as $\mathbf{X}_{car}^{(i)} = [\mathbf{x}_{car,1}^{(i)}, \mathbf{x}_{car,2}^{(i)}, \dots, \mathbf{x}_{car,C_{car}}^{(i)}]$. In this study, the aggregation function represents various fusion methods that will be introduced in the next section. q denotes neural networks that learn the latent representation (e.g., for CNNs, it is the feature maps) of the c th intermediate feature, where C_{car} and C_{mov} are the numbers of the intermediate features for cardiac sensing and movement sensing respectively.

3.3.3 Hybrid Fusion. With hybrid fusion, the fusion may occur at multiple stages/layers of the DL models [61, 83]. It's commonly understood that the DL model hierarchically encodes features at different levels, starting from low-level to higher-level features as the layers go deeper [31]. In this study, we would not cover all possible combinations of fusion architecture. Therefore, following previous work [61], we consider a simple scenario, which is firstly to fuse different input channel data belonging to the same modality (sharing a representation learning network) and then to fuse the high-level features from both modalities at the later stage. Formally, the hybrid-fusion strategy can be written as:

$$\hat{y}^{(i)} = \varphi(\text{Agg}(g_{mov}(\mathbf{X}_{mov}^{(i)}), g_{car}(\mathbf{X}_{car}^{(i)}))) \quad (4)$$

where φ is the classifier (e.g., fully connected neural networks) and g denotes the modality-specific networks (e.g., CNNs) that can learn representation from a specific modality such that the g_{mov} does not share network parameters with g_{car} . Agg is the aggregation function that can be implemented as concatenation[51], attention mechanism [79] and tensor-based method [22].

3.4 Fusion Method

Based on their complexity, fusion methods can be divided into three types: simple operations, attention-based methods and tensor-based methods. For feature vectors from different modalities, the concatenation and addition are two commonly used simple operations [85]. The attention mechanism is widely used for multimodal fusion. This usually refers to dynamically calculating a weight vector for each time step (or spatial position) and weighting a set of feature vectors [5, 18]. For tensor-based methods, bilinear pooling is a method of fusing two unimodal representations to a joint presentation by calculating their outer product. This method can capture the multiplicative interaction between all elements in two vectors [81].

For the early-stage fusion, we only adopted the concatenation as the only fusion method in this study. For the late-stage fusion, we selected two commonly used simple methods, which are concatenation and element-wise addition. Hybrid fusion provides aggregated representations for each modality, which facilitates flexible fusion methods. Apart from the simple operation methods, we also evaluated the attention mechanism and the tensor-based method. The choice of fusion method may be influenced by the application context.

3.4.1 Concatenation. For **early-stage fusion**, the concatenation method concatenates inputs of all modalities into one matrix, which can be denoted as:

$$\mathbf{K}_{early}^{(i)} = \text{Concatenate}(\mathbf{X}_{car}^{(i)}, \mathbf{X}_{mov}^{(i)}) \quad (5)$$

where $\mathbf{K}_{early}^{(i)} \in \mathbb{R}^{(C_{mov}+C_{car}) \times T}$. $\mathbf{X}_{car}^{(i)} \in \mathbb{R}^{C_{car} \times T}$ is the intermediate feature matrix of cardiac sensing, $\mathbf{X}_{mov}^{(i)} \in \mathbb{R}^{C_{mov} \times T}$ is the intermediate feature matrix of movement sensing, and the C_{car} represents the number of intermediate feature inputs of cardiac sensing.

For **late-stage fusion**, suppose we have the cardiac latent representation denoted as $\mathbf{X}'_{car,c} \in \mathbb{R}^{U \times L}$, which is learned from a neural network g via $\mathbf{X}'_{car,c} = g(\mathbf{x}_{car,c}^{(i)})$. The movement representation matrix is computed in the same way, which can be formally denoted as $\mathbf{X}'_{mov} = g(\mathbf{x}_{mov}^{(i)})$ where the $\mathbf{X}'_{mov} \in \mathbb{R}^{U \times L}$ is the latent representation of movement sensing. L is the temporal length and U is the representation's dimension. For example, in a convolutional neural network, U is the number of feature maps. At the late stage, as the feature maps of each input intermediate feature were kept separately, the concatenation operation concatenates these representations together, as follows:

$$\mathbf{K}_{late}^{(i)} = \text{Concatenate}(\mathbf{X}'_{mov,1}^{(i)}, \dots, \mathbf{X}'_{mov,C_{mov}}^{(i)}, \mathbf{X}'_{car,1}^{(i)}, \dots, \mathbf{X}'_{car,C_{car}}^{(i)}) \quad (6)$$

In this study, for the activity counts (handcraft feature) and cardiac features, the late-stage fusion's representation is denoted as $\mathbf{K}_{late}^{(i)} \in \mathbb{R}^{(C_{mov}+C_{car}) \times U \times L}$

For the **hybrid fusion**, the high-level representation of each modality is obtained from their own sub-network. The movement sensing representation is denoted as $\mathbf{X}''_{mov} = g_{mov}(\mathbf{X}_{mov}^{(i)})$ and the cardiac sensing is formally denoted as $\mathbf{X}''_{car} = g_{car}(\mathbf{X}_{car}^{(i)})$. The concatenation method for the hybrid fusion can be written as:

$$\mathbf{K}_{hybrid}^{(i)} = \text{Concatenate}(\mathbf{X}''_{car}, \mathbf{X}''_{mov}) \quad (7)$$

where $\mathbf{K}_{hybrid}^{(i)} \in \mathbb{R}^{2U \times L}$

3.4.2 Addition. The second simple operation is the element-wise addition denoted as \oplus . For the **late-stage fusion**, the addition operation is to integrate the high-level representation of each channel from each modality. The method is formally denoted:

$$\mathbf{Q}_{late}^{(i)} = \mathbf{X}'_{mov,1}^{(i)} \oplus, \dots, \mathbf{X}'_{mov,C_{mov}}^{(i)} \oplus \mathbf{X}'_{car,1}^{(i)}, \dots, \oplus \mathbf{X}'_{car,C_{car}}^{(i)} \quad (8)$$

where $\mathbf{Q}_{late}^{(i)} \in \mathbb{R}^{U \times L}$.

For the **hybrid fusion**, the addition method will aggregate the high-level representation of each modality. Formally, it can be denoted as:

$$\mathbf{Q}_{hybrid}^{(i)} = \mathbf{X}''_{car} \oplus \mathbf{X}''_{mov} \quad (9)$$

where $\mathbf{Q}_{hybrid}^{(i)} \in \mathbb{R}^{U \times T}$.

3.4.3 Attention Mechanism. Attention methods have been broadly adopted in multimodal fusion tasks. For example, in VQA tasks the method used is to fuse the visual representations with the language representation [79]. In our study we derived the attention model that could use the attention vectors to weight one modality based on the context of another modality. The meaning behind this is to filter the most significant information from a unimodal, which is jointly relevant for three-stage sleep classification. Therefore, we designed two attention fusion methods. The first one is Attention-on-Movement (Attention-on-Mov) and the second one is Attention-on-Cardiac (Attention-on-Car). Given the cardiac representation matrix $\mathbf{X}'_{car}^{(i)}$ and the movement representation

matrix $\mathbf{X}'_{mov}{}^{(i)}$, we firstly feed them through a single-layer neural network and then apply the softmax function to generate the attention distribution over the temporal dimension, which is denoted as:

$$\mathbf{H}_{att}^{(i)} = \tanh(\mathbf{W}_{car}\mathbf{X}'_{car}{}^{(i)} \oplus \mathbf{W}_{mov}\mathbf{X}'_{mov}{}^{(i)} + b_h) \quad (10)$$

$$\mathbf{P}_{att}^{(i)} = \text{softmax}(\mathbf{W}_{att}\mathbf{H}_{att}^{(i)} + b_{att}) \quad (11)$$

where $\mathbf{X}'_{mov}{}^{(i)} \in \mathbb{R}^{U \times L}$. Suppose we have linear transformation matrices that include $\mathbf{W}_{mov}, \mathbf{W}_{car} \in \mathbb{R}^{D \times U}$ and $\mathbf{W}_{att} \in \mathbb{R}^{L \times D}$, then $\mathbf{H}_{att}^{(i)} \in \mathbb{R}^{D \times L}$ and $\mathbf{P}_{att}^{(i)} \in \mathbb{R}^{L \times L}$, where D is the dimension of attention embedding space. The attention weight matrix is denoted as $\mathbf{P}_{att}^{(i)} = [\mathbf{p}_{att,1}^{(i)}, \dots, \mathbf{p}_{att,L}^{(i)}]$ and each temporal step has an attention vector $\mathbf{p}_{att,l}^{(i)}$, where $\sum \mathbf{p}_{att,l}^{(i)} = 1$. The subscript *att* stands for attention and l is the temporal step index.

We assume that applying attention weights on different modalities will have an impact on the results. Therefore, two scenarios were studied in this work. The first method is to weight cardiac sensing representation based on the attention distribution and concatenate them to build the joint feature representation matrix. It can be written as:

$$\mathbf{V}_{car}^{(i)} = \mathbf{X}'_{car}{}^{(i)}\mathbf{P}_{att}^{(i)} \quad (12)$$

$$\mathbf{K}_{car}^{(i)} = \text{Concatenate}(\mathbf{V}_{car}^{(i)}, \mathbf{X}'_{mov}{}^{(i)}) \quad (13)$$

We refer to this method as Attention-on-Car and $\mathbf{K}_{car}^{(i)} \in \mathbb{R}^{2U \times L}$

The second method is to weight the latent feature of movement sensing using the attention distribution, then concatenate them to build the joint representation matrix, which can be denoted as:

$$\mathbf{V}_{mov}^{(i)} = \mathbf{X}'_{mov}{}^{(i)}\mathbf{P}_{att}^{(i)} \quad (14)$$

$$\mathbf{K}_{mov}^{(i)} = \text{concatenate}(\mathbf{V}_{mov}^{(i)}, \mathbf{X}'_{car}{}^{(i)}) \quad (15)$$

We refer to this method as Attention-on-Mov and $\mathbf{K}_{mov}^{(i)} \in \mathbb{R}^{2U \times L}$ is the merged joint representation.

3.4.4 Bilinear Pooling Method. Bilinear pooling is a method to compute the matrices outer product that can facilitate multiplication interaction between all elements in both matrices. It's a method often used to fuse visual feature vectors with textual feature vectors to create a joint representation space, even though their distribution may vary dramatically[22, 45]. During the NREM sleep period, our cardiac system is co-modulated by peripheral and sympathetic neural systems. The heart rate is generally below the average of wake and REM sleep period and accompanied with tiny tremors in limb movement. We hypothesized that the bilinear model may be able to capture such tiny differences between REM and NREM sleep. Given its superior representation learning capacity, it has achieved remarkable performance in fine-grained image classification tasks [42]. Bilinear model calculates the outer product of two matrices. In this work, suppose we have two feature representation matrices $\mathbf{X}'_{car}{}^{(i)}$ and $\mathbf{X}'_{mov}{}^{(i)}$, and the bilinear representation can be written as:

$$\mathbf{k}_{bi}^{(i)} = \text{vec}(\mathbf{X}'_{car}{}^{(i)} \otimes \mathbf{X}'_{mov}{}^{(i)}) \quad (16)$$

The symbol of \otimes denotes Kronecker product of two matrices, and the *vec* denotes the matrix vectorization. After the vectorization, we then perform an element-wise signed square-root as denoted:

$$\mathbf{k}_{bi}^{(i)} \leftarrow \text{sign}(k_{bi}^{(i)})\sqrt{|k_{bi}^{(i)}|} \quad (17)$$

and then apply l_2 normalization on the vector $\mathbf{k}_{bi}^{(i)}$. We pass the normalized vector to a linear function that can reduce the feature dimensions before feeding them to the classifier.

4 EXPERIMENT DESIGN

In this section, we describe the experimental design of advanced multimodal fusion strategies and methods for the three-stage sleep classification using wearable devices. We firstly introduce two open-access datasets used in the study, including the data collection, data pre-processing and feature extraction. Secondly, we illustrate four backbone networks used with advanced multimodal fusion techniques. Finally, we list the evaluation metrics used in the study.

4.1 Dataset Description

4.1.1 Apple Watch Sleep Dataset. The first dataset used in our study is the Apple Watch Sleep Study¹, which is an open-access dataset collected at the University of Michigan between 2017 and 2019 [73, 74]. The dataset consists of 31 healthy subjects with no known sleep disorders or cardiovascular diseases and neurological or psychiatric impairment disorders [74]. All subjects wore Apple Watch (Apple Inc. series 2 and 3) and performed continuous recording for 7 to 14 days, and then joined the PSG study in the sleep laboratory on the last day [74]. During the PSG study, all subjects wore Apple Watch, which recorded heart rate and triaxial acceleration [74]. The acceleration and heart rate were measured by Apple Watch and recorded by a custom-developed watch application using the built-in functions of the iOS Watch kit and HealthKit by creating a “Workout Session” in app [73]. The PSG recordings were annotated according to the AASM rules [74].

The heart rate is measured by the PPG sensor of the Apple Watch and recorded as beats per minute (BPM), and a sample is returned from the Apple API every few seconds. The heart rate data is timestamped and the interval is between 2s and 5s. After the data cleaning process, we performed the feature engineering process on triaxial acceleration data; following [69, 74], we used activity counts as the movement feature. The final activity counts were added for each sleep epoch. Since the heart rate collected from Apple Watch is calculated in two to five seconds, we may treat them as a “pseudo” instantaneous heart rate (IHR). In each sleep epoch, we calculated the summary statistics of the heart rate data (called HR statistics or HRS for short), which includes the mean, standard deviation, minimum, maximum, skewness and kurtosis of the heart rate. Together with the activity counts, we constructed a seven-dimensional vector for each sleep epoch from these intermediate summary features and called it the Apple ACT-HRS feature set.

4.1.2 MESA Dataset. The Multi-Ethnic Study of Atherosclerosis (MESA) is a multi-site prospective study that includes 6,814 men and women. The ethnic groups include White, Black/African American, Hispanic, or Chinese, and the subjects are between the ages of 45 and 84 [14, 86]. The study was designed prospectively to evaluate risk factors of cardiovascular disease. The study had 2,237 participants enrolled in the sleep exam, which includes seven days of wrist-worn actigraphy; they underwent concurrent PSG for one night (wrist-worn actigraphy collected concurrently) [86]. The subjects who reported regular night-time use of nocturnal oxygen or positive airway pressure devices were excluded from the sleep exam [86]. The actigraphy recorded the activity counts in 1/30Hz and ECG is recorded at 100Hz.

We used the method and data processing protocol provided by the benchmark study [84] to synchronize PSG, ECG and activity counts of each subject. After the data pre-processing, 1,743 of 2,237 participants satisfied the data quality condition. Full details of the study setup, protocol and sampling rates are available in [14, 84, 86]. According to the feature set used in previous research [74, 84], we used the same features including activity counts and eight HRV features derived from the NN interval data in each sleep epoch [84]. The feature set consists of the Mean NNI, Standard Derivation of RR interval (SDNN), RR interval differences (SDSD), Very Low Frequency, Low Frequency, High Frequency Bands, Low Frequency to High Frequency Ratio and Total Power. These features have been investigated in several sleep physiology studies [12, 16, 20, 50, 71]. For each sleep epoch,

¹<https://physionet.org/content/sleep-accel/1.0.0/>

we constructed the intermediate feature vector based on eight HRV features and the activity counts (a scalar value per sleep epoch), which we named the MESA ACT-HRV feature set.

In addition, we converted the NN intervals into IHR data, and calculated the statistical features of IHR and combined activity counts as the second intermediate feature set. The purpose is to study the feature effects on the choices of fusion strategies and methods. We named it the MESA ACT-HRS feature set.

4.2 Evaluation Metrics

For performance evaluation, accuracy, Cohen’s κ , mean F1 and time deviation ([84]) were used. The time deviation that was used in the benchmark study [84] is denoted as $(TD_k = \frac{1}{N} \sum_{i=1}^N (Pred_c^i - GT_c^i))$. For a sleep stage c , the $Pred_c$ refers to the predicted minutes and GT_c refers to the ground truth sleep minutes. The superscript i represents the i th subject. The time deviation summarizes the mean bias of the total minutes of each sleep stage predicted by the classifier in the population. To understand the impact of individual differences in performance evaluation, this study adopted the subject-level evaluation. We calculated the metrics of each subject individually and obtained the mean value and 95% confidence interval of each metric for the population.

4.3 Experimental Procedure

Following previous work, we adopted a highly overlapping sliding window method with $S = 1$ to segment the input time-series data. In [84], the hyperparameter tuning results showed that the window length can impact the prediction performance. For convolutional neural networks, a longer window produced better results compared with a shorter window.

For each sleep epoch, we selected 50 adjacent (forward and backward) sleep epochs’ data to construct the inputs with a window length of 101. The details are shown in Figure 1. For the sliding window at the beginning and end of the recording, we filled these empty sleep epoch inputs with a value of -1. For the training, validation and testing, our experimental settings are as follows:

- **Apple Watch Sleep Dataset** Following the experiment setting of previous work [74], instead of using leave-one-subject-out-cross-validation, we adopted leave-two-subjects-out cross validation. Each fold had two subjects for testing, except for the last fold, which only contained one subject (total 31 subjects and 16 folds). In each fold, we then randomly split the subjects in training dataset into a validation dataset (20%) and a training dataset (80%). The validation set was used to select the best model for the test dataset.
- **MESA Sleep Dataset** The dataset contains 1743 valid sleep records of subjects. We employed the hold-out method to divide the entire dataset into a test set of 348 subjects (20%) and a training set of 1,395 subjects (80%) following the previous study [84]. The training set was further randomly split into a validation set (20%) and a training set (80%) [84]. Again, the validation set was used to select the best model for the test dataset.

All experiments conducted in this paper adopted the above setting for each dataset respectively.

In previous work [84], it was found that the performance improvement of three-stage sleep classification was more related to increasing the number of LSTM networks instead of increasing the number of CNN layers for the three-stage sleep classification task. Therefore, in this study, we focused on the design of CNN architecture. All experiments in this work adopted the Adam gradient update rule [37] with learning rate $\alpha = 10^{-4}$, $\beta_1 = 0.9$, and $\beta_2 = 0.99$. No early-stopping or weight decay was adopted in training processing. The batch size was set to 1024 except for the experiments containing the bilinear method which were set to 512. For the attention method, we set the attention embedding dimension to 256. For the bilinear method, we reduced the size of the feature dimension to 1024 using a linear layer. The training epoch corresponding to both datasets was set to 20.

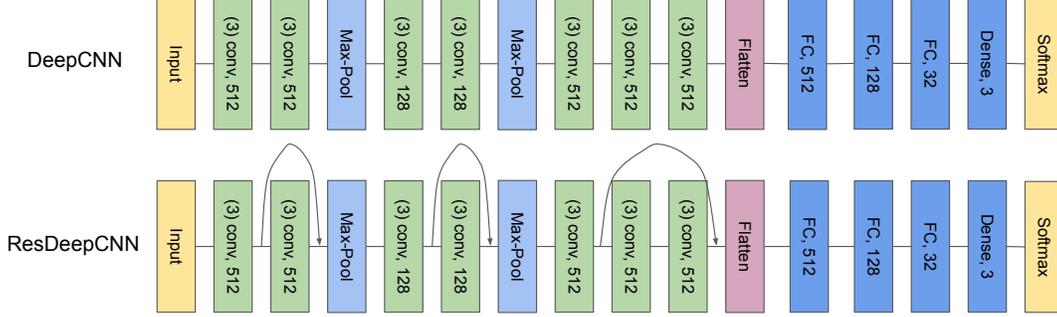


Fig. 2. Backbone network used in this study. The DeepCNN network was selected from the hyperparameter search results. We added the skip connection inside each convolutional block and we referred to it as ResDeepCNN. The stride and padding value were set to 1 for all convolutional (conv) layers, and the kernel size was set to 3. The kernel size and stride were set to 2 for all max pooling (Max-Pool) layers. The dropout was applied after each fully connected (FC) layer, and the dropout rate was set to 0.25.

4.4 Implementation Details

4.4.1 Hyperparameter Tuning and Backbone Networks. The fusion strategies and fusion methods may not benefit from a single layer convolutional neural network. To find a feasible backbone deep CNN that was capable of serving our study, inspired by [67], we designed a backbone network and conducted a hyperparameter search on 3-5 convolutional layer blocks (corresponding to 7-13 convolutional layers). From the hyperparameter tuning results, the network from the highest F1 validation score group was selected. We further gradually reduced the number of hidden units in fully connected layers and the experimental results showed slight improvements on model performance. We called this DeepCNN. To better understand the impact of modality fusion strategies and methods in different CNN architectures, inspired by [30, 53], we further added a skip connection in each convolutional block and called it ResDeepCNN, as the skip connection became an indispensable component in a variety of neural architectures that could boost representation learning. Figure 2 lists the details of two network structures. As our study focuses on the fusion strategy and methods, the backbone network was merely designed to conduct the feasible experiments. More details about the hyperparameter search can be seen in Appendix A

4.4.2 Backbone Network Setting. For the early-stage fusion and hybrid fusion, DeepCNN and ResDeepCNN were the main networks for the experiments. We slightly adapted the DeepCNN and the ResDeepCNN for the late-stage fusion experiments according to [78], which allowed each input channel to share the convolutional kernels but kept the feature representation separate. This means that, for a convolutional layer, the feature maps extracted from each input channel would not be fused with the feature maps of other channels. Instead, each input channel's feature maps would be fused before the classification module (fully-connected layers). For instance, for DeepCNN in the Apple Watch dataset, if the input was an intermediate feature matrix that contained cardiac and movement sensing and was denoted as $S_0^{(i)} \in \mathbb{R}^{7 \times 101}$ (one movement feature and six HRS features), the feature map function $\mathcal{F}_{l+1} : S_l^{(i)} \mapsto S_{l+1}^{(i)}$ was realized by a convolutional layer, where l denotes the l th convolutional layer. The output of the first convolutional layer was the feature map denoted as $S_1^{(i)} = \mathcal{F}_1(S_0^{(i)}) \in \mathbb{R}^{C_1 \times 7 \times 101}$, where C_1 was the number of feature maps of the first CNN layer. In this way, the intermediate feature of each input channel was kept separate.

5 RESULTS

This section empirically compares each combination of multimodal fusion strategies and methods based on two scenarios. The first scenario is the MESA dataset which contains multimodal data that can be extracted from research grade-wearable devices. The second scenario is the Apple Watch dataset derived from the consumer-grade smartwatches (Apple Watch Series 2 and 3) with the sleep stages annotated using the gold-standard PSG study.

We reported the performance in the order of three fusion strategies which included early-stage fusion, late-stage fusion, and hybrid fusion, and three fusion methods, including simple operation (concatenation and addition), attention mechanism, and bi-linear pooling method. We also investigated the effects of different window lengths (51 and 21), and the corresponding results can be seen in Appendix C.1. The experiments using raw accelerometer data and HR statistical features can be seen in Appendix B.2.

For consistency, all our fusion strategies and methods in each dataset were evaluated on the subject-level during the sleep recording period. Accuracy, Cohen’s κ , the mean F1 score and time deviation (minutes) were calculated based on the predictions during the sleep recording period. In the end, we compared the model parameter size and inference time for each strategy and method. These factors are important for model selection in the context of ubiquitous computing.

5.1 Apple Watch Dataset

5.1.1 Activity Counts and HRS Features. The first experiment was performed based on activity counts and HRS feature set (ACT-HRS) derived from the consumer wearables.

Table 1. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategies and methods with the Apple Watch dataset using ACT-HRS feature based on a window length of 101.

Fusion Specifics			Performance Metrics			Time Deviation(min.)		
Fusion Strategy	Network	Fusion Method	Accuracy(%)	Cohen’s κ	Mean F1(%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	72.3 \pm 2.5	40.0 \pm 5.8	59.0 \pm 3.6	-0.6 \pm 22.3	11.8 \pm 22.5	-11.2 \pm 6.9
	ResDeepCNN	Concatenation	76.0 \pm 2.4	45.7 \pm 6.1	63.4 \pm 3.5	12.3 \pm 17.2	-4.0 \pm 17.1	-8.2 \pm 7.3
Late-Stage Fusion	DeepCNN	Concatenation	76.2 \pm 2.7	47.0 \pm 7.1	63.7 \pm 4.5	8.4 \pm 16.5	1.3 \pm 19.4	-9.7 \pm 7.1
		Addition	78.2 \pm 2.1	49.9 \pm 6.9	65.2 \pm 3.8	11.4 \pm 10.4	-0.6 \pm 11.2	-10.8 \pm 6.6
	ResDeepCNN	Concatenation	75.5 \pm 2.9	47.0 \pm 6.5	63.4 \pm 4.2	10.4 \pm 18.5	-2.1 \pm 19.6	-8.3 \pm 7.0
		Addition	78.2 \pm 2.3	52.0 \pm 5.8	66.5 \pm 3.8	1.9 \pm 11.7	4.7 \pm 12.6	-6.5 \pm 7.3
Hybrid Fusion	DeepCNN	Concatenation	72.9 \pm 3.0	40.1 \pm 6.6	60.6 \pm 3.8	10.7 \pm 20.1	0.6 \pm 19.5	-11.4 \pm 6.4
		Addition	72.1 \pm 2.9	38.6 \pm 6.5	58.9 \pm 3.9	8.8 \pm 20.4	1.3 \pm 19.9	-10.2 \pm 7.3
		Attention-on-Mov	73.5 \pm 2.9	41.3 \pm 6.2	60.4 \pm 3.9	1.2 \pm 18.7	5.0 \pm 18.9	-6.2 \pm 7.7
		Attention-on-Car	71.1 \pm 3.4	37.4 \pm 6.7	58.1 \pm 4.1	7.5 \pm 24.3	0.7 \pm 23.1	-8.2 \pm 7.5
		Bilinear	69.5 \pm 3.5	29.5 \pm 7.5	53.0 \pm 4.2	1.2 \pm 25.7	10.0 \pm 25.0	-11.3 \pm 8.5
	ResDeepCNN	Concatenation	74.4 \pm 2.4	44.2 \pm 5.7	62.0 \pm 3.3	2.5 \pm 16.7	5.3 \pm 16.9	-7.8 \pm 7.0
		Addition	74.9 \pm 2.3	44.3 \pm 5.4	62.3 \pm 3.7	12.8 \pm 14.2	-7.6 \pm 15.6	-5.2 \pm 8.2
		Attention-on-Mov	75.2 \pm 2.6	45.0 \pm 5.6	63.1 \pm 3.4	10.4 \pm 19.4	-2.0 \pm 19.5	-8.4 \pm 7.2
		Attention-on-Car	72.2 \pm 3.0	41.5 \pm 6.8	59.7 \pm 3.8	-2.9 \pm 25.3	12.0 \pm 26.4	-9.1 \pm 6.7
		Bilinear	70.8 \pm 3.5	38.4 \pm 7.5	58.1 \pm 4.4	-2.8 \pm 22.9	6.4 \pm 24.8	-3.5 \pm 8.0

Table 1 lists the subject-level evaluation results of the Apple Watch dataset based on the window length of 101 during the sleep recording period. Since Apple Watch sampled the heart rate data with the unknown resolution and method, we only performed the experiments based on the ACT-HRS feature setting.

Overall, the ResDeepCNN achieved the highest mean F1 score of 66.5%, the Cohen’s κ of 52 and accuracy of 78.2% using the addition method in late-stage fusion. The same methods used on DeepCNN were higher than these in early-stage fusion too.

In the hybrid fusion strategy, using the Attention-on-Mov method achieved the highest scores irrespective of backbone networks. We observed the same pattern in the experiments using the window lengths 51 and 21. For window lengths 51 and 21, the highest performed models in each category were lower than the highest performed models using the window length of 101. We listed these results for window lengths 51 and 21 in Appendix C.1.

5.2 MESA Sleep Dataset Results

The second experiment was conducted on the MESA dataset. It was by far the largest sleep dataset that contained activity counts and instantaneous heart rate, which might be extracted from research-grade wearable devices. Again, we performed experiments on two different feature sets. The first feature set included activity counts and HRV features (ACT-HRV) [84], while the second feature set was ACT-HRS derived using the same feature extraction method in the Apple Watch dataset.

5.2.1 MESA Activity Counts and HRV Features. The reason for using the HRV features was that they had sleep physiological meaning. Table 2 shows the subject-level evaluation results based on the window length of 101. For the early-stage and late-stage fusion, the results of two backbone networks were comparable, which showed the skipping connections did not improve the classification performance.

Table 2. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategies and methods with the MESA test dataset using the ACT-HRV feature set based on a window length of 101.

Fusion Specifics			Performance Metrics			Time Deviation(min.)		
Fusion Strategy	Network	Fusion Method	Accuracy(%)	Cohen's κ	Mean F1(%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	78.6 \pm 0.9	62.8 \pm 1.8	71.1 \pm 1.3	27.1 \pm 6.9	-6.5 \pm 3.6	-20.7 \pm 6.4
	ResDeepCNN	Concatenation	78.0 \pm 1.1	60.2 \pm 2.0	71.1 \pm 1.4	54.1 \pm 7.0	1.2 \pm 3.8	-55.3 \pm 6.6
Late-Stage Fusion	DeepCNN	Concatenation	79.6 \pm 0.9	64.3 \pm 1.8	72.5 \pm 1.3	12.9 \pm 6.6	0.1 \pm 3.5	-13.0 \pm 6.2
		Addition	78.5 \pm 0.9	62.3 \pm 1.8	71.1 \pm 1.3	25.7 \pm 6.4	-6.8 \pm 3.3	-18.9 \pm 6.4
	ResDeepCNN	Concatenation	79.3 \pm 0.9	64.4 \pm 1.7	72.6 \pm 1.2	8.9 \pm 6.4	4.2 \pm 3.4	-13.1 \pm 6.1
		Addition	78.6 \pm 1.0	62.8 \pm 1.9	71.4 \pm 1.3	2.4 \pm 6.9	-3.0 \pm 3.5	0.6 \pm 6.7
Hybrid Fusion	DeepCNN	Concatenation	77.6 \pm 1.1	62.7 \pm 1.8	71.4 \pm 1.3	-12.7 \pm 7.3	7.2 \pm 3.9	5.5 \pm 7.0
		Addition	79.0 \pm 0.9	62.9 \pm 1.7	70.7 \pm 1.3	53.6 \pm 6.9	-15.0 \pm 3.3	-38.6 \pm 6.3
		Attention-on-Mov	79.0 \pm 1.0	63.9 \pm 1.8	72.1 \pm 1.3	2.3 \pm 7.0	-4.9 \pm 3.5	2.6 \pm 6.9
		Attention-on-Car	78.1 \pm 1.0	63.0 \pm 1.7	71.6 \pm 1.3	-2.1 \pm 6.8	6.5 \pm 4.0	-4.4 \pm 6.3
		Bilinear	75.7 \pm 0.9	58.6 \pm 1.8	68.8 \pm 1.2	3.7 \pm 6.8	16.6 \pm 4.0	-20.3 \pm 6.3
	ResDeepCNN	Concatenation	79.7 \pm 0.9	65.3 \pm 1.7	72.7 \pm 1.3	6.4 \pm 6.7	-7.7 \pm 3.4	1.3 \pm 6.7
		Addition	79.8 \pm 0.9	64.1 \pm 1.7	72.7 \pm 1.3	24.1 \pm 6.9	-9.7 \pm 3.2	-14.4 \pm 6.5
		Attention-on-Mov	79.6 \pm 1.0	65.5 \pm 1.8	73.3 \pm 1.3	-0.9 \pm 6.4	6.1 \pm 3.7	-5.1 \pm 6.3
		Attention-on-Car	78.5 \pm 1.0	62.7 \pm 1.7	70.5 \pm 1.3	37.6 \pm 7.0	-9.5 \pm 4.0	-28.1 \pm 6.2
		Bilinear	75.7 \pm 0.9	58.6 \pm 1.8	68.8 \pm 1.2	3.7 \pm 6.8	16.6 \pm 4.0	-20.3 \pm 6.3

For the hybrid fusion strategy, the ResDeepCNN achieved the highest accuracy, the Cohen's κ and the mean F1 score of 79.6%, 65.5 and 73.3%, respectively, using the Attention-on-Mov method. The results were statistically significant ($p < 0.05$) and higher than the models in the early-stage fusion. Those metrics were higher than the Attention-on-Car models too. Similar to the Apple Watch dataset, the performance of models based on window lengths 51 and 21 tend to be worse than experiments performed with window length 101. We list these results in Appendix C.1.

In terms of time deviation, DeepCNN achieved the optimal time deviation using the Attention-on-Mov method. The mean value of NREM sleep time deviation was 0.9, and the mean value of REM sleep time deviation was 6.1.

5.2.2 MESA Activity Counts and HRS Features. We derived the heart rate statistical features from the instantaneous heart rate (IHR) data in the MESA dataset. The purpose was to understand whether the type of intermediate feature would cause a difference in results.

The subject-level evaluation is shown in Table 3. In the early-stage fusion, similar to the ACT-HRV feature setting, the results of two backbone networks were comparable. The ResDeepCNN, using the Attention-on-Mov fusion method, achieved the highest accuracy, the Cohen’s κ , and the mean F1 score of 80.3%, 65.6, and 72.9%, respectively. However, the Attention-on-Mov model based on the ACT-HRS feature set highly overestimated the NREM sleep time and underestimated the wake minutes. Again, the models with window lengths 51 and 21 achieved lower performance than 101. Therefore, we listed these results in Appendix C.1

Table 3. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategies and methods with the MESA test dataset using the ACT-HRS feature set based on a window length of 101.

Fusion Specifics			Performance Metrics			Time Deviation(min.)		
Fusion Strategy	Network	Fusion Method	Accuracy(%)	Cohen’s κ	Mean F1(%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	78.0 \pm 1.0	63.4 \pm 1.8	72.0 \pm 1.2	2.7 \pm 7.0	15.6 \pm 4.0	-18.2 \pm 6.4
	ResDeepCNN	Concatenation	76.9 \pm 1.1	61.9 \pm 1.9	71.1 \pm 1.3	-36.7 \pm 7.8	12.1 \pm 4.2	24.5 \pm 7.5
Late-Stage Fusion	DeepCNN	Concatenation	79.1 \pm 1.0	64.7 \pm 1.7	72.8 \pm 1.3	0.6 \pm 6.9	7.5 \pm 3.7	-8.1 \pm 6.4
	ResDeepCNN	Addition	78.1 \pm 0.9	62.8 \pm 1.6	70.6 \pm 1.2	23.1 \pm 6.6	-4.5 \pm 3.7	-18.5 \pm 6.3
	DeepCNN	Concatenation	77.8 \pm 1.0	62.5 \pm 1.7	71.0 \pm 1.3	-1.1 \pm 7.3	-0.7 \pm 3.5	1.8 \pm 6.8
	ResDeepCNN	Addition	77.7 \pm 1.0	62.6 \pm 1.7	70.7 \pm 1.2	24.3 \pm 7.0	3.7 \pm 3.9	-28.0 \pm 6.4
Hybrid Fusion	DeepCNN	Concatenation	78.2 \pm 0.9	64.4 \pm 1.7	70.2 \pm 1.2	18.5 \pm 7.1	-14.6 \pm 3.3	-3.9 \pm 6.5
		Addition	78.1 \pm 0.9	62.2 \pm 1.8	71.2 \pm 1.2	14.7 \pm 7.4	1.4 \pm 3.6	-16.1 \pm 6.8
		Attention-on-Mov	79.2 \pm 0.9	63.8 \pm 1.8	71.8 \pm 1.3	28.1 \pm 7.3	-4.3 \pm 3.5	-23.9 \pm 6.5
		Attention-on-Car	76.6 \pm 1.0	61.4 \pm 1.7	70.4 \pm 1.2	-7.4 \pm 7.9	17.7 \pm 4.7	-10.3 \pm 6.7
	ResDeepCNN	Bilinear	75.6 \pm 0.9	58.0 \pm 1.8	67.7 \pm 1.2	6.3 \pm 7.6	-8.1 \pm 3.7	1.8 \pm 7.1
		Concatenation	79.4 \pm 1.0	64.4 \pm 1.7	72.7 \pm 1.2	31.7 \pm 6.9	4.9 \pm 3.6	-36.6 \pm 6.3
		Addition	78.9 \pm 0.9	63.6 \pm 1.8	72.2 \pm 1.2	24.6 \pm 7.0	4.9 \pm 3.5	-29.5 \pm 6.5
		Attention-on-Mov	80.3 \pm 0.9	65.6 \pm 1.7	72.9 \pm 1.3	35.5 \pm 6.9	0.8 \pm 3.6	-36.3 \pm 6.2
	Attention-on-Car	79.3 \pm 0.9	62.8 \pm 1.7	71.1 \pm 1.2	29.1 \pm 7.1	-2.4 \pm 3.7	-26.7 \pm 6.3	
	Bilinear	74.1 \pm 0.9	56.8 \pm 1.7	66.9 \pm 1.2	6.2 \pm 7.1	11.7 \pm 4.2	-17.9 \pm 6.5	

5.3 Inference Efficiency

In the mobile computing scenario of three-sleep stage classification, the model based on the deep learning architecture may require sufficient computing resources. This may be a challenge for many inexpensive or low-end smartwatches and smartphones. Table 4 shows the model parameter size and inference time of each combination of fusion strategies and methods. In addition, we counted the number of trainable parameters and calculated the time required for forward propagation (running on CPU). All experiments were conducted using Pytorch 1.6 and the hardware platform consisted of 8 cores AMD-7 3700X with 4.4GHz and 64GB DDR4 memory. We independently ran each model 10 times on the Apple Watch dataset. Each time, we inferred 500 samples (sleep epochs) using the Pytorch profiling module to calculate the statistical summary of the inference time.

Overall, the models using the addition method in late-stage fusion, hybrid fusion, and concatenation in early-stage fusion had the least model parameters. As a result, the models in early-stage fusion achieved the shortest inference time. The addition method in the late-stage fusion had the same number of parameters as the models in early-stage fusion, but the inference time was increased by 7-10 times. This was because the late-stage fusion calculated the feature maps of each input channel separately and fused them before the classifier module (fully connected layers). Consequently, the feature matrix extracted by the convolutional module was a 3D tensor (e.g., the number of input feature dimensions \times number of feature maps \times temporal steps) in the late-stage fusion.

Table 4. The number of model parameters and inference time of each combination of fusion strategies and methods evaluated in millions of parameters and milliseconds respectively with the *Apple Watch* dataset, using the *ACT-HRS* feature sets based on a window length of 101.

Fusion Strategy	Network	Fusion Method	Total Parameters (M)	Inference Time (ms per sample)
Early-Stage	DeepCNN	Concatenation	9.44	3.52±0.08
	ResDeepCNN	Concatenation	9.44	3.54±0.08
Late-Stage Fusion	DeepCNN	Concatenation	48.75	22.9±0.17
		Addition	9.43	32.25±5.53
	ResDeepCNN	Concatenation	48.75	31.16±5.06
		Addition	9.43	22.11±0.25
Hybrid Fusion	DeepCNN	Concatenation	18.80	7.13±0.12
		Addition	12.24	7.01±0.12
		Attention-on-Act	19.07	7.32±0.12
		Bilinear	274.65	10.09±0.11
	ResDeepCNN	Concatenation	18.80	7.02±0.16
		Addition	12.24	7.0±0.14
		Attention-on-Act	19.07	7.27±0.17
		Bilinear	274.65	10.22±0.17

In contrast, the early-stage fusion generated a 2D tensor (e.g., number of feature maps \times temporal steps). The convolution operation required additional time to calculate the feature maps of each input channel.

The bilinear model had the largest model parameters. Most model parameters belonged to the feature representation module, which contained a fully connected layer to reduce the dimension of feature representation at an order of two magnitudes. Since the calculation speed of the fully connected (FC) layer was much faster than the convolutional layer, the inference time did not increase as much as the model parameter size.

6 EXPLORATION OF USING GRAD-CAM ON SLEEP SENSING DATA

One of the major drawbacks to the consumer sleep monitoring devices was inaccurate results with an unknown decision making process, which would harm the user’s confidence and trust in the devices [63]. Explanation of the decision-making process lies at the heart of a responsible research in applied machine learning [27]. Before building confidence and trust, the first step is communication and explanation [41]. Unlike multimedia data, the time-series data requires the user to associate a meaning to the values in the temporal dimension [27].

To investigate in what kind of visualization of the decision-making process of multimodal fusion can be understood by humans, the gradient class activation map (Grad-CAM) [66] was adopted to visualize the important areas that matter to a specific class prediction from a qualitative perspective. In sleep physiology, the HRV features such as HF and LF have been proven to be different based on different sleep stages, e.g., NREM sleep is associated with a lower overall HRV, and REM sleep is accompanied by increased variability [72]. Wake is associated most often with changes in activity counts. Based on these phenomena, we selected subjects from the MESA test dataset with F1 scores greater than 90%. To obtain clear graphics, a post-processing method was used on the Grad-CAM output, which simplified the time-series data by setting the heat map value to 1 if the CAM value was greater than a threshold of 0.8, otherwise it was set to 0.2.

Not all sleep epochs generated have consistent patterns, as the neural network is known to be able to learn background information that is relevant to the classification [1, 27]. Then we chose sleep epochs with good quality ². We then presented the visualization results to sleep experts, and the feedback showed that sleep technologists tend to use fewer PSG channels to reduce the overload of irrelevant information during sleep stage annotation.

²a) stay in a sleep stage for at least 5 minutes. b) The fluctuation patterns of the highlighted areas should be similar to that described in the studies of sleep physiology.

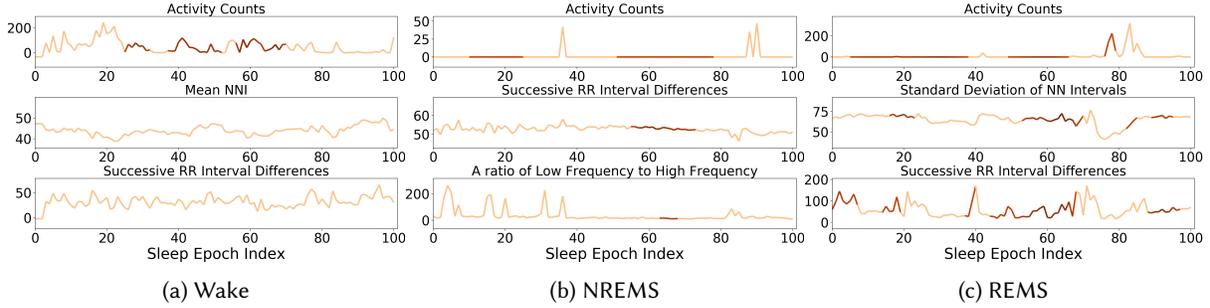


Fig. 3. The Grad-CAM plot of three selected examples from MESA dataset (ACT-HRV feature) using ResDeepCNN (Addition) in the late-stage fusion. Each row is the activation map for the input clinical features. To obtain a clear graph, the highlighted areas are the activation values over the threshold of 0.8 and the light color areas represents the activation values below the threshold of 0.8

To reduce the input channels, filtering them by feature importance scores calculated on the CAM value is a feasible method. We randomly selected 1k correctly predicted samples for each sleep stage from the test set and calculated their CAM values under the window length of 101. For each sleep epoch, we counted the number of time steps if their CAM value was greater than 0.8 and summed and normalized them as occurrences for each intermediate feature to obtain their relative importance score. Afterwards, we then calculated the mean value of all samples for each intermediate feature per sleep stage. Figure 4 shows the feature importance for each sleep stage.

We retained the top three channels of each sleep stage to simplify the visualization as shown in Figure 3. To test whether this visualization is useful and can be understood by humans, we designed a game system³ that could conducted the exploratory study with users. The game system serves the purpose of engaging user to read and understand these visualization. The study consists of two phases, which investigate the accuracy of sleep stage classification by humans based solely on the input signals in the cases of Non-CAM and CAM visualization, respectively. In each phase, we encoded a continuous period of sleep data (intermediate feature data and hypnogram) for each sleep stage into videos to speed up the training process. In each phase, users would first watch the training videos; then they would be asked to recognize nine randomly selected sleep epochs that did not belong to the training videos. At the end of the test, users will be informed of their sleep stage classification accuracy.

There are not established rules for sleep stage classification using movement and cardiac sensing data. As a pilot study, we recruited 25 individuals from Amazon’s Mechanical Turk. The task took, on average, 20-45 minutes to complete and participants were compensated USD 7.00. All procedures received ethical approval from the University’s ethical review board and the Research Ethics Committees (RECs). In total, we received 25 answers. Figure 5 (a) shows the classification accuracy of CAM-assisted sleep stage classification is higher than the Non-CAM sleep stage classification. We further analyzed the results in detail by sleep stages in Figure 5 (c). As can be seen, CAM-assisted visualizations improved human recognition accuracy on all sleep stages. To test whether the system can improve user’s understanding of the visualization, we also designed a question a five point Likert scale. The results showed that the majority of participants either *Strongly Agree* or *Agree* that the machine-assisted visualization helped them to understand the difference between each sleep stage.

³<https://gradcamvisual1.azurewebsites.net/>

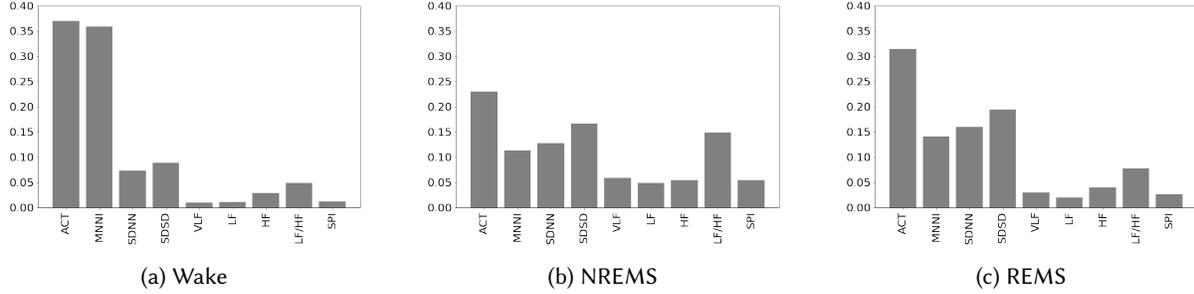


Fig. 4. The mean of total class activation value for each sleep stage from MESA dataset (ACT-HRV feature) using ResDeepCNN (Addition) in the late-stage fusion. ACT : Activity Counts, MNNI : Mean NNI, SDNN: Standard Deviation of NNI, SDSO : Successive RR Interval Differences, VLF: Very-Low-Frequency Band, LF: Low-Frequency Band, HF: High-Frequency Band, LF/HF: The ratio of Low Frequency to High Frequency, SPI: The Signal Power Intensity

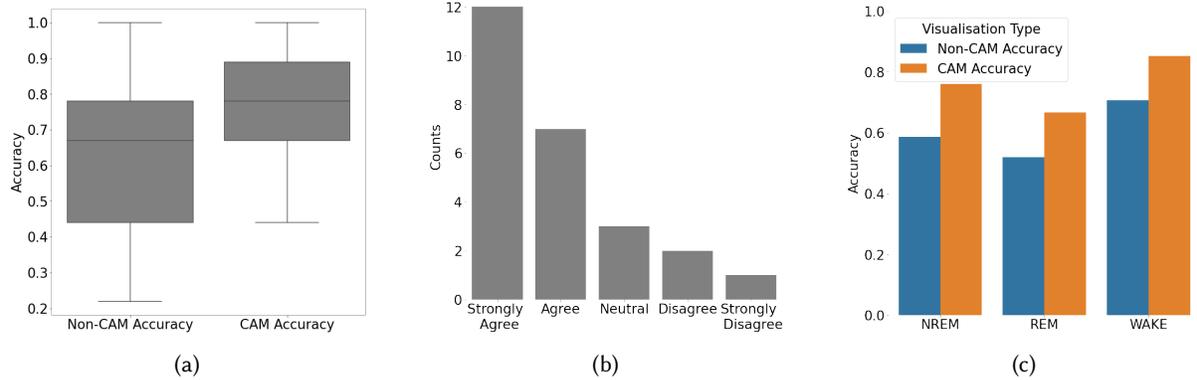


Fig. 5. (a) A user study of three-stage classification accuracy calculated per participant-wise for Non-CAM and CAM assisted visualization. (b) The answer of *The machine-assisted visualization helped me to understand the difference between each sleep stage*. (c) The breakdown of classification accuracy calculated for each sleep stage.

7 DISCUSSION

7.1 Simple Fusion Method and Fusion Strategy

7.1.1 Concatenation and Addition. All three fusion strategies used the concatenation method. In the early-stage fusion, the backbone network fused the latent features from each modality at every convolutional layer. All models studied in this paper surpassed previous studies, except for the models using the bilinear method. With the Apple Watch dataset, the skip connection numerically improved the performance for models in early-stage fusion, but the performance with the MESA dataset decreased. A possible explanation for this might be that simply adding a skip connection may not benefit the model prediction performance in early-stage fusion when the training data is sufficient.

In the late-stage fusion, the concatenation method numerically improved the prediction performance of all metrics for ResDeepCNN compared with the early-stage fusion. In terms of the concatenation method, parameter size and reasoning time increased by five times and six times, respectively, but the performance did not increase by that much. The addition method in late-stage fusion achieved the highest performance with the Apple Watch

dataset. This may indicate the benefit of keeping latent features separate, and fusing them at a higher level might produce better results. Another possible explanation for this is the increase in network parameters.

For the hybrid fusion, we first fused cardiac intermediate features at the early stage and fused them with the movement sensing representation at the late stage. With the Apple Watch dataset, the simple operation method produced the same or better results compared with early-stage fusion. The methods in this category increased model parameters and inference time, but the classification performance hardly improved. The addition method aggregated the modal representation at the later stage, and similarly, it only obtained comparable results. A similar pattern was observed with the MESA dataset. These findings suggest that the effectiveness of simple operation methods may be affected by fusion strategies.

7.2 Complex Fusion Method and Fusion Strategy

Bilinear pooling turned out to be the weakest fusion method in the hybrid fusion strategy. This is a particularly interesting result, because the tensor-based methods showed improvements in the multimodal fusion literature such as with the task of visual question answering [81]. It is possible that these results could be due to the failure to use the CNN network to learn about a post-bilinear latent feature. The model parameters were too large to be suitable for mobile computing scenarios. The cost of exploring the potential solutions exceeded the benefits.

With the attention mechanism, with the MESA dataset using the ACT-HRV feature setting, the mean F1 and Cohen's κ of ResDeepCNN using the Attention-on-Mov method were statistically higher than the those in early-stage fusion. With the Apple Watch dataset, the same method can also produce comparable results to the highest performing method in the late-stage fusion, and the inference speed is three times faster. Moreover, the time deviation value was balanced in the prediction of each sleep stage.

Compared to the Attention-on-Car method, we observed a similar pattern with the MESA dataset, that is, applying attention weights to the latent features of movement produced higher results. It can therefore be assumed that the attention method improved the network's ability to learn better representations that can benefit three-stage sleep classification by adjusting the weights of movement sensing representations, as it was difficult to discern REM sleep and NREM sleep using movement sensing alone.

7.3 Model Selection

The model parameters, model architecture, model inference time and model performance were the key considerations for the model selection in ubiquitous computing. In addition to these factors, the time deviation should also be considered. It reflected the bias of the model prediction for each sleep stage. With imbalanced sleep data sets, biased models may constantly overestimate the duration of certain sleep stages and may lead to unreasonable health decisions. The mean and standard error of time deviation should be as close to zero as possible. In terms of inference time, predicting sleep data for a whole night, using the late fusion strategy was 6.2 times slower than when using the early-stage fusion strategy. As the model calculates the feature maps of each input channel individually using the convolution method, the time consumption of this method was related to the number of input intermediate features. In addition, for the design of the fusion strategy, we should also consider the ratio of module parameters to inference time. For instance, the parameters of the bilinear model were 77 times larger than the early-stage fusion models, but the speed was only three times slower. because most of the model parameters belonged to a fully connected layer in the bilinear module, which reduced the feature dimension of the feature matrix (the results of the matrix outer product).

From the results of Appendix B.2, the use of raw accelerometer data generated comparable results when compared to the use of handcraft features but with increased model parameters and inference time. It was considered a sub-optimal solution for long-term sleep stage monitoring.

In summary, if the movement and cardiac sensing data can be transmitted to the smartphone, the ResDeepCNN with the addition method in the late-stage fusion may be a feasible model for everyday use. Since it achieved the highest F1 score and with a balanced time deviation on each sleep stage. In the scenarios with limited computing resources such as smartwatches, using the early-stage fusion and the ResDeepCNN model may be a practical choice. The cost of using the late-stage fusion has increased by 9 times in inference time.

7.4 Cross Dataset Comparison

Based on activity counts and HRS features, the highest performing model for the MESA dataset achieved higher performance than the highest such model for the Apple Watch dataset. Not only were the classification metrics of the MESA dataset higher than these of the Apple Watch dataset, but the standard error on the MESA dataset was smaller. It seems possible that these differences were due to two reasons. The first reason was that the Apple Watch dataset contained far fewer subjects compared with the MESA dataset. The second possible reason was the differences in data acquisition equipment and data pre-processing methods. The HR sensing module of Apple Watch dynamically calculated the HR data within two–five seconds, whereas the cardiac sensing used in the MESA dataset is IHR. The higher resolution IHR data might provide more discriminant information in order to discern three sleep stages.

7.5 Exploratory Research of Visualization

One of the objectives of this work is to better understand in what way the decision-making process of neural networks using multimodal fusion techniques on sleep stage classification can be understood by humans. Based on Grad-CAM scores, a simplified visualization method was adopted in this study and an exploratory study with users was performed. Our visualization tool can highlight both the key temporal signal segments and the most discriminant feature channels, and by providing important clues/patterns for users it can serve as an assistant tool in understanding different sleep stage signals.

Compared with highlighting the temporal steps, the channel dimension reduction retained the minimum number of discriminant channels for sleep stage classification, which showed a combined reduction that could improve user understanding. Based on the repeated patterns of activity count and HRV features during a continuous sleep period, the results demonstrated that reducing information overload could improve human understanding on three-stage sleep recognition performance. The visualization increased users' understanding in terms of the neural network decision making process to some extent.

Wake recognition accuracy was higher than in the other two sleep stages, which indicating the wrist movement, is obvious to classify wake. This result is consistent with the known capability of actigraphy can be used to distinguish between wake and sleep. The highlighted patterns that appeared in the continuous sleep stage agreed with previous sleep physiology findings to some extent.

The modeling process has window bias and Grad-CAM modeling bias. For example, this study adopted a window length of 101 (50.5 minutes). The highlighted areas will move backwards as the window moving forwards. The network is capable of locating signatures in a window that are meaningful for the current sleep stage recognition. This is very different to the annotation process using high-resolution (over 100Hz) PSG data. An interesting question for future work is to investigate whether these patterns have physiological meaning.

Many existing interpretation and visualization techniques have been developed for visual data, yet it is unclear whether these methods are suitable for explaining sleep time-series data. This is a pilot study to investigate whether Grad-CAM applied to sleep time-series wearable data may be useful to humans. It is one of the mainstream methods used in visual and text data but is subject to the network structure. For instance, it is difficult to highlight the important areas in channels in early-stage fusion models without substantial changes to the method. On the other hand, it is difficult for humans to understand the highlighted patterns in the time dimension.

We also observed the interesting results, e.g., 3 out of 25 users experienced negative impacts on their understandings. In addition, the questionnaire data also highlighted that not everyone agreed that the system helped them to understand the difference in the patterns between sleep stages. A possible explanation for this might be that the visualization may not be understood by every person, or the training and testing process may be problematic. Future studies may consider conducting experiments with more detailed personalized questions.

The visualization was designed as a pilot study to understand the decision-making process of multimodal fusion for sleep stage classification. So, we did not investigate other mainstream interpretation methods, nor did we conduct large-scale user research. Future work may consider investigating the other interpretation methods such as SHAP, Anchor, etc. or even create a special algorithm for time-series data to facilitate intuitive understanding by humans.

7.6 Comparison with Previous Work and Implications

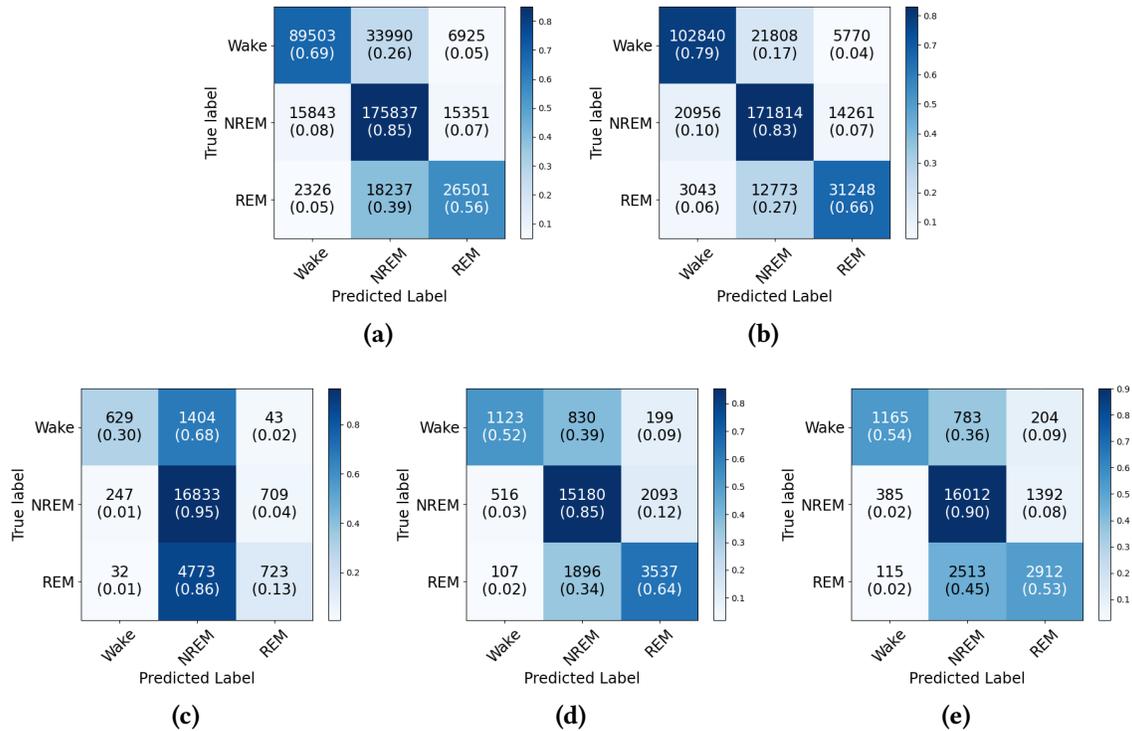


Fig. 6. **(a)** The CNN(101) and early-stage fusion based on the MESA (ACT-HRV) dataset used in [84]. **(b)** Hybrid fusion using ResDeepCNN (Attention-on-Act) based on the MESA (ACT-HRV) dataset **(c)** Walch et al. using multiple layer perceptron based on activity counts, HR and circadian time [74] **(d)** Late-stage fusion using ResDeepCNN (Addition) based on the Apple Watch Dataset **(e)** Late-stage fusion using ResDeepMixCNN (Concatenation) using raw accelerometer data and HRS based on the Apple Watch dataset.

Table 5. Three-stage sleep classification prediction results compared with previous work evaluated at subject level (mean \pm standard error at 95% confidence interval) during the recording period.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Dataset and Feature Set	Fusion Stage	Model	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
MESA (ACT-HRV)	Early-stage Fusion	CNN (101) (Zhai, 2020)	76.0 \pm 1.0	58.6 \pm 1.9	68.1 \pm 1.3	14.9 \pm 6.7	-0.5 \pm 4.3	-14.4 \pm 5.8
	Hybrid Fusion	ResDeepCNN (Attention-on-Mov)	79.8 \pm 0.9	65.5 \pm 1.8	73.3 \pm 1.3	-0.9 \pm 6.4	6.1 \pm 3.7	-5.1 \pm 6.3
Apple Watch (Activity Counts, HR, Time)	Early-stage Fusion	MLP (Walch, 2019)	72.1 \pm 2.4	23.7 \pm 4.4	47.8 \pm 3.6	84.2 \pm 17.2	-65.4 \pm 15.0	-18.8 \pm 6.1
Apple Watch (Activity Counts, HRS)	Late-stage Fusion	ResDeepCNN Addition	78.2 \pm 2.3	52.0 \pm 5.8	66.5 \pm 3.8	1.9 \pm 11.7	4.7 \pm 12.6	-6.5 \pm 7.3

Table 5 and Figure 6 show the results compared with previous works. We have observed that the use of multimodal fusion strategies and fusion methods can improve model prediction performance. With the three-stage sleep classification dataset, the class imbalance issue causes the classifier to be biased towards the majority class, which is NREM sleep.

To compare with the previous work, we conducted ten runs of the model with the highest mean F1 and the baseline model for each dataset, respectively. Each run used a different random number seed. We performed a t -test on the accuracy, Cohen's κ and mean F1 score. Compared with the previous work[84] with the MESA dataset, the accuracy ($p < 0.001$), Cohen's κ ($p < 0.001$) and mean F1 score ($p < 0.001$) improved statistically significantly with the MESA dataset using the ACT-HRV feature set. With the Apple Watch dataset, the accuracy ($p < 0.001$), mean F1 score ($p < 0.001$) and Cohen's κ ($p < 0.001$) were also statistically higher than previous work [74]. These improvements suggest that the proper multimodal fusion strategy and method can improve the robustness of the model, which is a step towards automated three-stage sleep classification. The findings reported here suggest that reasonable performance may be achieved using the movement and cardiac features derived from consumer/research-grade wearable devices.

8 CONCLUSION

Using actigraphy to monitoring sleep-wake has existed for many decades. But, for sleep stage classification, we have relied on the PSG study, which is an expensive, burdensome, laboratory sleep monitoring method. This limits many research advances in sleep and health. In recent years, more and more new products using ubiquitous computing technology have been passed FDA clearance, such as the Apple Watch irregular heart rate detection function [56]. The achievements of these wearables provide important instrumental tools for study of longitudinal sleep and health. The core contribution of this work lies in our systematic study on how to better integrate multi-modal data to monitor three-stage sleep that may use consumer/research grade wearables. Through our study, the performance of several new models exceeded the previous benchmark studies significantly. We have provided a new multimodal fusion benchmark for the ubiquitous computing community. This has provided the potential for the use of consumer wearables to study the sleep health of large-scale populations in the future. One of the motivations for this work was to respond to previous research that called for more accurate and transparent sleep stage sensing algorithms on consumer wearable devices [74]. We hope this work will encourage more researchers, consumers, and application developers to use consumer/research grade wearables to study and understand sleep and health.

ACKNOWLEDGMENTS

This research is funded through the EPSRC Centre for Doctoral Training in Digital Civics (EP/L016176/1). The Multi-Ethnic Study of Atherosclerosis (MESA) Sleep Ancillary study was funded by NIH-NHLBI Association

of Sleep Disorders with Cardiovascular Health Across Ethnic Groups (RO1 HL098433). MESA is supported by NHLBI funded contracts HHSN268201500003I, N01-HC-95159, N01-HC-95160, N01-HC-95161, N01-HC-95162, N01-HC-95163, N01-HC-95164, N01-HC-95165, N01-HC-95166, N01-HC-95167, N01-HC-95168 and N01-HC-95169 from the National Heart, Lung, and Blood Institute, and by cooperative agreements UL1-TR-000040, UL1-TR-001079, and UL1-TR-001420 funded by NCATS. The National Sleep Research Resource was supported by the National Heart, Lung, and Blood Institute (R24 HL114473, 75N92019R002). Thanks Duan Haoran, Ouyang Zizhou, Shao Shuai, Dr. Kirstie Anderson, Becky Zhu for their suggestions on model selection, experimental setup, and visualization, and sharing sleep knowledge with me.

REFERENCES

- [1] Amina Adadi and Mohammed Berrada. 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access* 6 (9 2018), 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- [2] Triantafyllos Afouras, Joon Son Chung, Andrew Senior, Oriol Vinyals, and Andrew Zisserman. 2018. Deep Audio-visual Speech Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2018). <https://doi.org/10.1109/TPAMI.2018.2889052>
- [3] Mehmet Y. Agargun and Rosalind Cartwright. 2003. REM sleep, dream variables and suicidality in depressed patients. *Psychiatry Research* 119, 1-2 (7 2003), 33–39. [https://doi.org/10.1016/S0165-1781\(03\)00111-2](https://doi.org/10.1016/S0165-1781(03)00111-2)
- [4] Sonia Ancoli-Israel, Roger Cole, Cathy Alessi, Mark Chambers, William Moorcroft, and Charles P. Pollak. 2003. The role of actigraphy in the study of sleep and circadian rhythms. , 342–392 pages. <https://doi.org/10.1093/sleep/26.3.342>
- [5] Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. International Conference on Learning Representations, ICLR.
- [6] Yang Bai, Yu Guan, and Wan Fai Ng. 2020. Fatigue assessment using ECG and actigraphy sensors. In *Proceedings - International Symposium on Wearable Computers, ISWC*. Association for Computing Machinery, 12–16. <https://doi.org/10.1145/3410531.3414308>
- [7] Eti Ben Simon, Aubrey Rossi, Allison G. Harvey, and Matthew P. Walker. 2020. Overanxious and underslept. *Nature Human Behaviour* 4, 1 (1 2020), 100–110. <https://doi.org/10.1038/s41562-019-0754-8>
- [8] Richard Berry. 2012. *Fundamentals of Sleep Medicine*. Elsevier Inc. <https://doi.org/10.1016/C2009-0-38997-7>
- [9] Richard B. Berry, Rita Brooks, Charlene Gamaldo, Susan M. Harding, Robin M. Lloyd, Stuart F. Quan, Matthew T. Troester, and Bradley V. Vaughn. 2017. AASM scoring manual updates for 2017 (version 2.4). , 665–666 pages. <https://doi.org/10.5664/jcsm.6576>
- [10] Luciana Besedovsky, Tanja Lange, and Jan Born. 2012. Sleep and immune function. , 121–137 pages. <https://doi.org/10.1007/s00424-011-1044-0>
- [11] Martin G. Bleichner, Micha Lundbeck, Matthias Selisky, Falk Minow, Manuela Jäger, Reiner Emkes, Stefan Debener, and Maarten De Vos. 2015. Exploring miniaturized EEG electrodes for brain-computer interfaces. An EEG you do not see? *Physiological Reports* 3, 4 (2015). <https://doi.org/10.14814/phy2.12362>
- [12] Philippe Boudreau, Wei Hsien Yeh, Guy A. Dumont, and Diane B. Boivin. 2013. Circadian variation of heart rate variability across sleep stages. *Sleep* 36, 12 (12 2013), 1919–1928. <https://doi.org/10.5665/sleep.3230>
- [13] Ramona Cabiddu, Sergio Cerutti, Geoffrey Viardot, Sandra Werner, and Anna M. Bianchi. 2012. Modulation of the sympatho-vagal balance during sleep: Frequency domain study of heart rate variability and respiration. *Frontiers in Physiology* 3 MAR (2012). <https://doi.org/10.3389/fphys.2012.00045>
- [14] Xiaoli Chen, Rui Wang, Phyllis Zee, Pamela L. Lutsey, Sogol Javaheri, Carmela Alcántara, Chandra L. Jackson, Michelle A. Williams, and Susan Redline. 2015. Racial/ethnic differences in sleep disturbances: The Multi-Ethnic Study of Atherosclerosis (MESA). *Sleep* 38, 6 (6 2015), 877–888. <https://doi.org/10.5665/sleep.4732>
- [15] Yuanying Chen, Wei Dong, Yi Gao, Xue Liu, and Tao Gu. 2017. Rapid. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 3 (9 2017), 1–27. <https://doi.org/10.1145/3130906>
- [16] Florian Chouchou and Martin Desseilles. 2014. Heart rate variability: A tool to explore the sleeping brain? , 402 pages. <https://doi.org/10.3389/fnins.2014.00402>
- [17] Maarten De Vos and Stefan Debener. 2014. Mobile eeg: Towards brain activity monitoring during natural action and cognition. , 2 pages. <https://doi.org/10.1016/j.ijpsycho.2013.10.008>
- [18] Haoran Duan, Shidong Wang, and Yu Guan. 2020. SOFA-Net: Second-Order and First-order Attention Network for Crowd Counting. (8 2020). <https://arxiv.org/abs/2008.03723v1>
- [19] Tasali E, Leproult R, Ehrmann DA, and Van Cauter E. 2008. Slow-wave sleep and the risk of type 2 diabetes in humans. *Proceedings of the National Academy of Sciences of the United States of America* 105, 3 (1 2008), 1044–1049. <https://doi.org/10.1073/PNAS.0706446105>
- [20] Pedro Fonseca, Niek Den Teuling, Xi Long, and Ronald M. Aarts. 2017. Cardiorespiratory Sleep Stage Detection Using Conditional Random Fields. *IEEE Journal of Biomedical and Health Informatics* 21, 4 (2017), 956–966. <https://doi.org/10.1109/JBHI.2016.2550104>

- [21] Pedro Fonseca, Niek Den Teuling, Xi Long, and Ronald M. Aarts. 2018. A comparison of probabilistic classifiers for sleep stage classification. *Physiological Measurement* 39, 5 (2018). <https://doi.org/10.1088/1361-6579/aabbc2>
- [22] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. 2016. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP 2016 - Conference on Empirical Methods in Natural Language Processing, Proceedings*. Association for Computational Linguistics (ACL), 457–468. <https://doi.org/10.18653/v1/d16-1044>
- [23] Nina E. Fultz, Giorgio Bonmassar, Kawin Setsompop, Robert A. Stickgold, Bruce R. Rosen, Jonathan R. Polimeni, and Laura D. Lewis. 2019. Coupled electrophysiological, hemodynamic, and cerebrospinal fluid oscillations in human sleep. *Science* 366, 6465 (11 2019), 628–631. <https://doi.org/10.1126/science.aax5440>
- [24] Laurent Giovannardi, Omer T. Inan, Richard M. Wiard, Mozziyar Etemadi, and Gregory T.A. Kovacs. 2011. Ballistocardiography - A method worth revisiting. In *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBS*, Vol. 2011. NIH Public Access, 4279–4282. <https://doi.org/10.1109/IEMBS.2011.6091062>
- [25] Yu Guan, Chang Tsun Li, and Fabio Roli. 2015. On Reducing the Effect of Covariate Factors in Gait Recognition: A Classifier Ensemble Method. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37, 7 (7 2015), 1521–1528. <https://doi.org/10.1109/TPAMI.2014.2366766>
- [26] Yu Guan and Thomas Plötz. 2017. Ensembles of Deep LSTM Learners for Activity Recognition using Wearables. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 2 (3 2017), 1–28. <https://doi.org/10.1145/3090076>
- [27] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *Comput. Surveys* 51, 5 (8 2018). <https://doi.org/10.1145/3236009>
- [28] Junichiro Hayano, Emi Yuda, and Yutaka Yoshida. 2017. Sleep stage classification by combination of actigraphic and heart rate signals. *2017 IEEE International Conference on Consumer Electronics - Taiwan, ICCE-TW 2017* (2017), 387–388. <https://doi.org/10.1109/ICCE-China.2017.7991158>
- [29] Caner Hazirbas, Lingni Ma, Csaba Domokos, and Daniel Cremers. 2017. FuseNet: Incorporating depth into semantic segmentation via fusion-based CNN architecture. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Vol. 10111 LNCS. Springer Verlag, 213–228. https://doi.org/10.1007/978-3-319-54181-5_14
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Identity mappings in deep residual networks. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, Vol. 9908 LNCS. Springer Verlag, 630–645. https://doi.org/10.1007/978-3-319-46493-0_38
- [31] Jeff Heaton. 2018. Ian Goodfellow, Yoshua Bengio, and Aaron Courville: Deep learning. *Genetic Programming and Evolvable Machines* 19, 1-2 (6 2018), 305–307. <https://doi.org/10.1007/s10710-017-9314-z>
- [32] Chen-Yu Hsu, Aayush Ahuja, Shichao Yue, Rumen Hristov, Zachary Kabelac, and Dina Katabi. 2017. Zero-Effort In-Home Sleep and Insomnia Monitoring using Radio Signals. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 1, 3 (9 2017), 1–18. <https://doi.org/10.1145/3130924>
- [33] IBER and C. 2007. The AASM Manual for the Scoring of Sleep and Associated Events : Rules. *Terminology and Technical Specification* (2007). <https://ci.nii.ac.jp/naid/10024500923>
- [34] Ashish Kapoor and Rosalind W. Picard. 2005. Multimodal affect recognition in learning environments. In *Proceedings of the 13th ACM International Conference on Multimedia, MM 2005*. ACM Press, New York, New York, USA, 677–682. <https://doi.org/10.1145/1101149.1101300>
- [35] Vishesh K. Kapur, Dennis H. Auckley, Susmita Chowdhuri, David C. Kuhlmann, Reena Mehra, Kannan Ramar, and Christopher G. Harrod. 2017. Clinical Practice Guideline for Diagnostic Testing for Adult Obstructive Sleep Apnea: An American Academy of Sleep Medicine Clinical Practice Guideline. *Journal of Clinical Sleep Medicine : JCSM : Official Publication of the American Academy of Sleep Medicine* 13, 3 (2017), 479. <https://doi.org/10.5664/JCSM.6506>
- [36] Seema Khosla, Maryann C. Deak, Dominic Gault, Cathy A. Goldstein, Dennis Hwang, Younghoon Kwon, Daniel O’Hearn, Sharon Schutte-Rodin, Michael Yurcheshen, Ilene M. Rosen, Douglas B. Kirsch, Ronald D. Chervin, Kelly A. Carden, Kannan Ramar, R. Nisha Aurora, David A. Kristo, Raman K. Malhotra, Jennifer L. Martin, Eric J. Olson, Carol L. Rosen, and James A. Rowley. 2018. Consumer sleep technology: An American academy of sleep medicine position statement. *Journal of Clinical Sleep Medicine* 14, 5 (5 2018), 877–880. <https://doi.org/10.5664/jcsm.7128>
- [37] Diederik P. Kingma and Jimmy Lei Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. International Conference on Learning Representations, ICLR. <https://arxiv.org/abs/1412.6980v9>
- [38] Michelle Kohansieh and Amgad N. Makaryus. 2015. Sleep Deficiency and Deprivation Leading to Cardiovascular Disease. <https://doi.org/10.1155/2015/615681>
- [39] David J. Kupfer and F. Gordon Foster. 1972. INTERVAL BETWEEN ONSET OF SLEEP AND RAPID-EYE-MOVEMENT SLEEP AS AN INDICATOR OF DEPRESSION. *The Lancet* 300, 7779 (9 1972), 684–686. [https://doi.org/10.1016/S0140-6736\(72\)92090-9](https://doi.org/10.1016/S0140-6736(72)92090-9)
- [40] Palagini L, Baglioni C, Ciapparelli A, Gemignani A, and Riemann D. 2013. REM sleep dysregulation in depression: state of the art. *Sleep medicine reviews* 17, 5 (10 2013), 377–390. <https://doi.org/10.1016/J.SMRV.2012.11.001>

- [41] Brian Y. Lim, Anind K. Dey, and Daniel Avrahami. 2009. Why and why not explanations improve the intelligibility of context-aware intelligent systems. *Conference on Human Factors in Computing Systems - Proceedings (2009)*, 2119–2128. <https://doi.org/10.1145/1518701.1519023>
- [42] Tsung Yu Lin, Aruni Roychowdhury, and Subhransu Maji. 2015. Bilinear CNN models for fine-grained visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, Vol. 2015 Inter. 1449–1457. <https://doi.org/10.1109/ICCV.2015.170>
- [43] Shengzhong Liu, Shuochao Yao, Jinyang Li, Dongxin Liu, Tianshi Wang, Huajie Shao, and Tarek Abdelzaher. 2020. GlobalFusion: A global attentional deep learning framework for multisensor information fusion. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 1 (3 2020), 1–27. <https://doi.org/10.1145/3380999>
- [44] David Looney, Preben Kidmose, Cheolsoo Park, Michael Ungstrup, Mike Rank, Karin Rosenkranz, and Danilo Mandic. 2012. The in-the-ear recording concept: User-centered and wearable brain monitoring. *IEEE Pulse* 3, 6 (2012), 32–42. <https://doi.org/10.1109/MPUL.2012.2216717>
- [45] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical question-image co-attention for visual question answering. In *Advances in Neural Information Processing Systems*. Neural information processing systems foundation, 289–297.
- [46] A. Malliani, M. Pagani, F. Lombardi, and S. Cerutti. 1991. Cardiovascular neural regulation explored in the frequency domain. *Circulation* 84, 2 (1991), 482–492. <https://doi.org/10.1161/01.CIR.84.2.482>
- [47] M. Méndez, A. M. Bianchi, O. Villantieri, and S. Cerutti. 2006. Time-varying analysis of the heart rate variability during REM and non REM sleep stages. In *Annual International Conference of the IEEE Engineering in Medicine and Biology - Proceedings*. 3576–3579. <https://doi.org/10.1109/IEMBS.2006.260067>
- [48] Kaare B. Mikkelsen, Yousef R. Tabar, Simon L. Kappel, Christian B. Christensen, Hans O. Toft, Martin C. Hemmsen, Mike L. Rank, Marit Otto, and Preben Kidmose. 2019. Accurate whole-night sleep monitoring with dry-contact ear-EEG. *Scientific Reports* 9, 1 (12 2019), 1–12. <https://doi.org/10.1038/s41598-019-53115-3>
- [49] Kaare B. Mikkelsen, David Bové Villadsen, Marit Otto, and Preben Kidmose. 2017. Automatic sleep staging using ear-EEG. *BioMedical Engineering Online* 16, 1 (9 2017), 111. <https://doi.org/10.1186/s12938-017-0400-5>
- [50] Nicola Montano, Alberto Porta, Chiara Cogliati, Giorgio Costantino, Eleonora Tobaldini, Karina Rabello Casali, and Ferdinando Iellamo. 2009. Heart rate variability explored in the frequency domain: A tool to investigate the link between heart and behavior. , 71–80 pages. <https://doi.org/10.1016/j.neubiorev.2008.07.006>
- [51] A Monti, C. Medigue, H. Nedelcoux, and P. Escourrou. 2002. Autonomic control of the cardiovascular system during sleep in normal subjects. *European Journal of Applied Physiology* 87, 2 (6 2002), 174–181. <https://doi.org/10.1007/s00421-002-0597-1>
- [52] Natalia Neverova, Christian Wolf, Graham Taylor, and Florian Nebout. 2016. ModDrop: Adaptive multi-modal gesture recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 38, 8 (8 2016), 1692–1706. <https://doi.org/10.1109/TPAMI.2015.2461544>
- [53] A. Emin Orhan and Xaq Pitkow. 2017. Skip Connections Eliminate Singularities. *arXiv* (1 2017). <http://arxiv.org/abs/1701.09175>
- [54] Joao Palotti, Raghvendra Mall, Michael Aupetit, Michael Rueschman, Meghna Singh, Aarti Sathyanarayana, Shahrads Taheri, and Luis Fernandez-Luque. 2019. Benchmark on a large cohort for sleep-wake classification with machine learning techniques. *npj Digital Medicine* 2, 1 (12 2019), 1–9. <https://doi.org/10.1038/s41746-019-0126-9>
- [55] Kwang Suk Park, Su Hwan Hwang, Da Woon Jung, Hee Nam Yoon, and Won Kyu Lee. 2014. Ballistocardiography for noninvasive sleep structure estimation. In *2014 36th Annual International Conference of the IEEE Engineering in Medicine and Biology Society, EMBC 2014*. Institute of Electrical and Electronics Engineers Inc., 5184–5187. <https://doi.org/10.1109/EMBC.2014.6944793>
- [56] Marco V. Perez, Kenneth W. Mahaffey, Haley Hedlin, John S. Rumsfeld, Ariadna Garcia, Todd Ferris, Vidhya Balasubramanian, Andrea M. Russo, Amol Rajmane, Lauren Cheung, Grace Hung, Justin Lee, Peter Kowey, Nisha Talati, Divya Nag, Santosh E. Gummidipundi, Alexis Beatty, Mellanie True Hills, Sumbul Desai, Christopher B. Granger, Manisha Desai, and Mintu P. Turakhia. 2019. Large-Scale Assessment of a Smartwatch to Identify Atrial Fibrillation. *New England Journal of Medicine* 381, 20 (11 2019), 1909–1917. <https://doi.org/10.1056/NEJMoa1901183>
- [57] Ignacio Perez-Pozuelo, Bing Zhai, Joao Palotti, Raghvendra Mall, Michaël Aupetit, Juan M. Garcia-Gomez, Shahrads Taheri, Yu Guan, and Luis Fernandez-Luque. 2020. The future of sleep health: a data-driven revolution in sleep science and medicine. , 15 pages. <https://doi.org/10.1038/s41746-020-0244-4>
- [58] Brandon Peters. 2021. Overnight Sleep Study: Uses, Procedure, Results. <https://www.verywellhealth.com/what-to-expect-in-a-sleep-study-3015121>
- [59] Huy Phan, Fernando Andreotti, Navin Cooray, Oliver Y. Chen, and Maarten De Vos. 2019. SeqSleepNet: End-to-End Hierarchical Recurrent Neural Network for Sequence-to-Sequence Automatic Sleep Staging. *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 27, 3 (2019), 400–410. <https://doi.org/10.1109/TNSRE.2019.2896659>
- [60] Huy Phan, Oliver Y. Chen, Minh C. Tran, Philipp Koch, Alfred Mertins, and Maarten De Vos. 2021. XSleepNet: Multi-View Sequential Model for Automatic Sleep Staging. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021). <https://doi.org/10.1109/TPAMI.2021.3070057>
- [61] Valentin Radu, Catherine Tong, Sourav Bhattacharya, Nicholas D Lane, Cecilia Mascolo, Mahesh K. Marina, and Fahim Kawsar. 2018. Multimodal Deep Learning for Activity and Context Recognition. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous*

- Technologies* 1, 4 (2018), 1–27. <https://doi.org/10.1145/3161174>
- [62] A. Rattani, D. R. Kisku, M. Bicego, and M. Tistarelli. 2007. Feature level fusion of face and fingerprint biometrics. In *IEEE Conference on Biometrics: Theory, Applications and Systems, BTAS'07*. <https://doi.org/10.1109/BTAS.2007.4401919>
- [63] Ruth Ravichandran, Sang Wha Sien, Shwetak N. Patel, Julie A. Kientz, and Laura R. Pina. 2017. Making sense of sleep sensors: How sleep sensing technologies support and undermine sleep health. *Conference on Human Factors in Computing Systems - Proceedings 2017-May* (5 2017), 6864–6875. <https://doi.org/10.1145/3025453.3025557>
- [64] Reutrakul S and Van Cauter E. 2018. Sleep influences on obesity, insulin resistance, and risk of type 2 diabetes. *Metabolism: clinical and experimental* 84 (7 2018), 56–66. <https://doi.org/10.1016/J.METABOL.2018.02.010>
- [65] Jonathan Schwartz and Thomas Roth. 2009. Neurophysiology of Sleep and Wakefulness: Basic Science and Clinical Implications. *Current Neuropharmacology* 6, 4 (2 2009), 367–378. <https://doi.org/10.2174/157015908787386050>
- [66] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2020. Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. *International Journal of Computer Vision* 128, 2 (2020), 336–359. <https://doi.org/10.1007/s11263-019-01228-7>
- [67] Karen Simonyan and Andrew Zisserman. 2015. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*. International Conference on Learning Representations, ICLR. <http://www.robots.ox.ac.uk/>
- [68] Ambra Stefani and Birgit Högl. 2019. Sleep in Parkinson's disease. *Neuropsychopharmacology* 2019 45:1 45, 1 (6 2019), 121–128. <https://doi.org/10.1038/s41386-019-0448-y>
- [69] Bart H.W. Te Lindert and Eus J.W. Van Someren. 2013. Sleep estimates using microelectromechanical systems (MEMS). *Sleep* 36, 5 (2013), 781–789. <https://doi.org/10.5665/sleep.2648>
- [70] Jing Xian Teo, Sonia Davila, Chengxi Yang, An An Hii, Chee Jian Pua, Jonathan Yap, Swee Yaw Tan, Anders Sahlén, Calvin Woon-Loong Chin, Bin Tean Teh, Steven G. Rozen, Stuart Alexander Cook, Khung Keong Yeo, Patrick Tan, and Weng Khong Lim. 2019. Digital phenotyping by consumer wearables identifies sleep-associated markers of cardiovascular disease risk and biological aging. *Communications Biology* 2019 2:1 2, 1 (10 2019), 1–10. <https://doi.org/10.1038/s42003-019-0605-1>
- [71] Eleonora Tobaldini, Lino Nobili, Silvia Strada, Karina R. Casali, Alberto Braghiroli, and Nicola Montano. 2013. Heart rate variability in normal and pathological sleep. *Frontiers in Physiology* 4 (10 2013), 1–11. <https://doi.org/10.3389/fphys.2013.00294>
- [72] Emilio Vanoli, Philip B. Adamson, Ba-Lin, Gian D. Pinna, Ralph Lazzara, and William C. Orr. 1995. Heart Rate Variability During Specific Sleep Stages. *Circulation* 91, 7 (4 1995), 1918–1922. <https://doi.org/10.1161/01.CIR.91.7.1918>
- [73] Olivia Walch. 2019. Motion and heart rate from a wrist-worn wearable and labeled sleep from polysomnography v1.0.0. <https://physionet.org/content/sleep-accel/1.0.0/>
- [74] Olivia Walch, Yitong Huang, Daniel Forger, and Cathy Goldstein. 2019. Sleep stage prediction with raw acceleration and photoplethysmography heart rate data derived from a consumer wearable device. *Sleep* 42, 12 (12 2019). <https://doi.org/10.1093/sleep/zsz180>
- [75] Yi-Qun Wang, Rui Li, Meng-Qi Zhang, Ze Zhang, Wei-Min Qu, and Zhi-Li Huang. 2015. The Neurobiological Mechanisms and Treatments of REM Sleep Disturbances in Depression. *Current Neuropharmacology* 13, 4 (9 2015), 543. <https://doi.org/10.2174/1570159X13666150310002540>
- [76] Zeyu Wang, Xiongfei Li, Haoran Duan, Yanchi Su, Xiaoli Zhang, and Xinjiang Guan. 2021. Medical image fusion based on convolutional neural networks and non-subsampled contourlet transform. *Expert Systems with Applications* 171 (6 2021), 114574. <https://doi.org/10.1016/J.ESWA.2021.114574>
- [77] Zeyu Wang, Xiongfei Li, Haoran Duan, Xiaoli Zhang, and Hancheng Wang. 2019. Multifocus image fusion using convolutional neural networks in the discrete wavelet transform domain. *Multimedia Tools and Applications* 2019 78:24 78, 24 (8 2019), 34483–34512. <https://doi.org/10.1007/S11042-019-08070-6>
- [78] Jian Bo Yang, Minh Nhut Nguyen, Phyo Phyo San, Xiao Li Li, and Shonali Krishnaswamy. 2015. Deep convolutional neural networks on multichannel time series for human activity recognition. In *IJCAI International Joint Conference on Artificial Intelligence*, Vol. 2015-Janua. 3995–4001.
- [79] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. 2016. Stacked attention networks for image question answering. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Vol. 2016-Decem. IEEE Computer Society, 21–29. <https://doi.org/10.1109/CVPR.2016.10>
- [80] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Vol. 2016-Decem. IEEE Computer Society, 4651–4659. <https://doi.org/10.1109/CVPR.2016.503>
- [81] Zhou Yu, Jun Yu, Jianping Fan, and Dacheng Tao. 2017. Multi-modal Factorized Bilinear Pooling with Co-attention Learning for Visual Question Answering. In *Proceedings of the IEEE International Conference on Computer Vision*, Vol. 2017-Octob. Institute of Electrical and Electronics Engineers Inc., 1839–1848. <https://doi.org/10.1109/ICCV.2017.202>
- [82] Shichao Yue, Yuzhe Yang, Hao Wang, Hariharan Rahul, and Dina Katabi. 2020. BodyCompass: Monitoring Sleep Posture with Wireless Signals. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 2 (6 2020), 1–25. <https://doi.org/10.1145/3628888>

//doi.org/10.1145/3397311

[83] Xin Zeng, Qiang Guo, Haoran Duan, and Yunpeng Wu. 2021. Multi-level features extraction network with gating mechanism for crowd counting. *IET Image Processing* (2021). <https://doi.org/10.1049/IPR2.12304>

[84] Bing Zhai, Ignacio Perez-Pozuelo, Emma A.D. Clifton, Joao Palotti, and Yu Guan. 2020. Making Sense of Sleep: Multimodal Sleep Stage Classification in a Large, Diverse Population Using Movement and Cardiac Sensing. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies* 4, 2 (6 2020), 1–33. <https://doi.org/10.1145/3397325>

[85] Chao Zhang, Zichao Yang, Xiaodong He, and Li Deng. 2020. Multimodal Intelligence: Representation Learning, Information Fusion, and Applications. *IEEE Journal on Selected Topics in Signal Processing* 14, 3 (3 2020), 478–493. <https://doi.org/10.1109/JSTSP.2020.2987728>

[86] Guo Qiang Zhang, Licong Cui, Remo Mueller, Shiqiang Tao, Matthew Kim, Michael Rueschman, Sara Mariani, Daniel Mobley, and Susan Redline. 2018. The National Sleep Research Resource: Towards a sleep data commons. *Journal of the American Medical Informatics Association* 25, 10 (10 2018), 1351–1358. <https://doi.org/10.1093/jamia/ocy064>

A HYPERPARAMETERS TUNING AND RESULTS

The hyperparameter tuning was performed based on the designed backbone network from three to five convolutional blocks (7-13 convolutional layers). The first two blocks consisted of two convolutional layers. The third, fourth and fifth convolutional blocks consisted of three convolutional layers. The hyperparameter search aimed to reduce the search space and maintain suitable temporal lengths of the latent features. The hyperparameter tuning only focused on the number of kernels for each convolutional block. The convolutional layer kernel length has been investigated in the previous study [84]. We set the kernel length of all convolutional layers to 3.

The number of hidden units in the fully connected layers was all set to the same value during the hyperparameter tuning process to reduce the search space. Furthermore, we performed the hyperparameter tuning based on the MESA dataset - the largest dataset containing the cardiac and activity data to date. Therefore, we expected the hyperparameter tuning could discover robust backbone networks for this study.

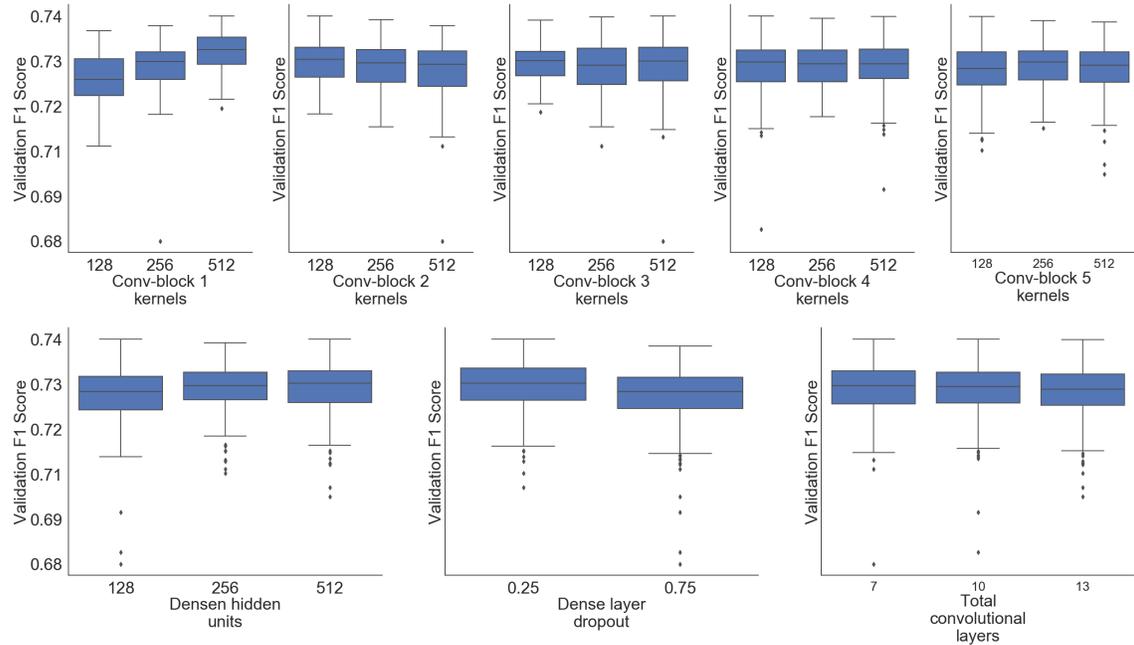


Fig. 7. DeepCNN backbone network hyper-parameters tuning results

Table 6. Hyper-parameters tuning for backbone networks

Block/Group	Layer	Kernel size	Number of Kernels	Padding	Stride	Hidden Units	Drop out Rate
Convolutional Block 1	Convolutional Layer 1 & 2	3	128, 256, 512	1	1		
	Maxpooling	2			2		
Convolutional Block 2	Convolutional Layer 3 & 4	3	128, 256, 512	1	1		
	Maxpooling	2			2		
Convolutional Block 3	Convolutional Layer 5, 6 & 7	3	128, 256, 512	1	1		
	Maxpooling	2			2		
Convolutional Block 4	Convolutional Layer 8, 9 & 10	3	128, 256, 512	1	1		
	Maxpooling	2			2		
Convolutional Block 4	Convolutional Layer 11, 12 & 13	3	128, 256, 512	1	1		
	Maxpooling	2			2		
FC Block	Fully Connected Layer 1, 2					128, 256, 512	
	Drop Out						0.25, 0.75

B HEART RATE STATISTIC FEATURES COMBINED WITH DEEP MOVEMENT FEATURES

In our study, we also tested whether using the raw accelerometer data could achieve better results. Therefore, we designed two feasible CNNs that could extract the compatible deep features fused with the HR statistic features. The rationale behind the network design was to produce a compatible representation of the HR intermediate feature. To match up the dimension of latent feature, we firstly reduced the accelerometer data sampling rate from 50Hz/20Hz to 1Hz. We then designed two CNNs to bridge the sampling gap between movement and cardiac sensing features, one for the early stage fusion and another for the late-stage and hybrid fusion. Each consisted of six convolutional layers (two convolutional blocks) to extract the deep movement feature used for fusion study. For the early-stage fusion, the network was called AccCNN-1. The hybrid and late-stage fusion used the same network to extract the latent representations, and we called it AccCNN-2. We referred to the entire network as DeepMixCNN and ResDeepMixCNN, respectively. The network structure and the experiment setting details can be seen in Figure 8 and Figure 9. We adopted the leave-two-subjects-out cross validation experimental setting on Apple Watch dataset. The training, validation and testing process used the same settings as the main content. However, we did not conduct the hyperparameter search together with the backbone network. Therefore, the network designed in the study merely served as feasible networks for the study, yet it might not be the best performing CNN. We focused on the fusion techniques rather than the contribution of network structure.

The MESA dataset contained the activity counts sampled at 1/30 Hz, which technically was not the raw data. In addition, cardiac sensing was acquired via the PSG equipment, which may be difficult to wear everyday. Therefore, the HRV features derived from the RR intervals were most likely to be available from the commercial wearable devices (e.g., photoplethysmogram data), so we did not conduct the experiments on the raw PSG data. The details of this experiment are listed in Table 7

B.1 Raw Accelerometer Data and HRS Features

. The highest performed model was ResDeepMixCNN in late-stage fusion, using the concatenation method. Its accuracy, the Cohen's κ score and the mean F1 reached 79.1 %, 51.4 and 66.7 % respectively. Thus, the results were comparable to the handcraft features.

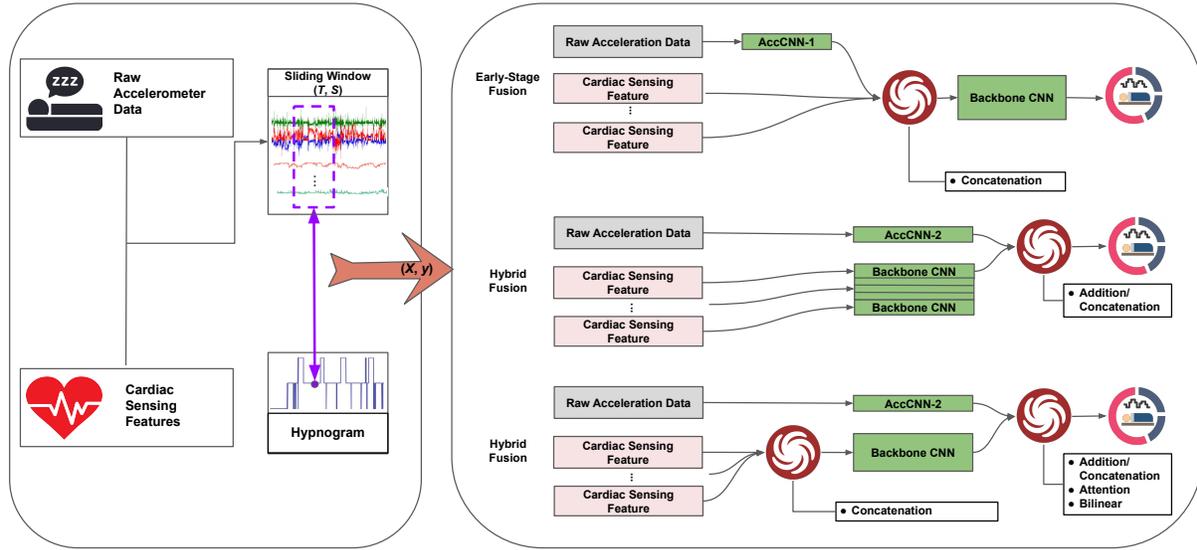


Fig. 8. An overview of the three-stage sleep classification system using the raw accelerometer data with HR statistics features. The raw accelerometer data and HR statistic features were extracted for each sleep epoch (30s). The sliding window method divides the sleep data into multiple segments with window length T and stride S . In this experiment, we have $T = 101$, and $S = 1$. We firstly use the AccCNN to learn deep features then fuse them with HR statistic features. The hypnogram represents the stages of sleep over time. Two fusion strategies and four fusion methods were studied.

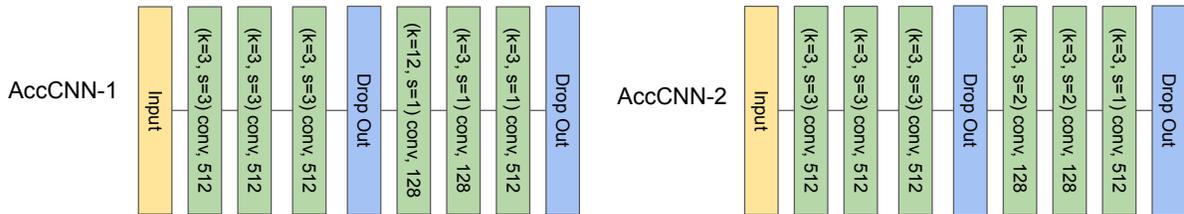


Fig. 9. An overview of the two subnets used to extract the deep features from the raw accelerometer data.

B.2 Comparison of Raw Data and Intermediate Features

We compared the performance difference between using the raw accelerometer data and using the clinical/handcraft features based on the window length of 101. The ResDeepMixCNN has achieved the comparable performance on the Apple Watch dataset in terms of accuracy, Cohen’s κ and mean F1, using the concatenation method in the late-stage fusion. The confusion matrices shown in Figure 6 demonstrated the model prediction using raw accelerometer data is biased to NREM sleep. Three reasons might cause the increased bias. The first reason may be the Apple Watch dataset has class imbalance issue. The second reason may be the modality bias of the raw accelerometer data because the wrist movement may not reflect the sleep stage (mainly NREM and REM) that much. The third reason may be caused by the lacking of hyperparameter search on the network.

Table 7. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) using raw accelerometer data and HRS features based on DeepMixCNN and ResDeepMixCNN with the Apple Watch Dataset for each combination of fusion strategy and method. The experiments were performed using the same experimental setting as in the main content and evaluated at the subject level during recording period based on window length of 101.

Fusion Specifics			Performance Metrics			Deployment Metrics	
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Model Size (M)	Inference Time (ms)
Early-Stage	DeepMixCNN	Concatenation	74.8 \pm 2.7	34.0 \pm 5.9	57.2 \pm 3.5	11.9	18.86 \pm 1.19
	ResDeepMixCNN	Concatenation	79.8 \pm 3.1	48.9 \pm 7.1	64.5 \pm 4.5	11.9	16.45 \pm 0.43
Late-Stage Fusion	DeepMixCNN	Concatenation	79.2 \pm 2.8	48.9 \pm 7.6	66.0 \pm 4.6	50.8	37.03 \pm 0.34
		Addition	76.7 \pm 2.1	38.7 \pm 5.9	58.0 \pm 3.7	11.5	32.75 \pm 0.19
	ResDeepMixCNN	Concatenation	79.1 \pm 3.3	51.4 \pm 8.0	66.7 \pm 4.7	50.8	31.19 \pm 0.72
		Addition	78.9 \pm 2.8	48.4 \pm 6.4	63.9 \pm 4.0	11.5	33.11 \pm 0.33
Hybrid	DeepMixCNN	Concatenation	77.3 \pm 3.3	43.9 \pm 7.3	63.6 \pm 4.3	18.0	15.94 \pm 0.22
		Addition	75.8 \pm 2.6	36.7 \pm 7.3	59.0 \pm 4.0	11.5	16.28 \pm 0.57
		Attention-on-Mov	75.8 \pm 3.1	39.4 \pm 7.9	60.8 \pm 4.8	18.3	15.7 \pm 0.18
		Attention-on-Car	71.9 \pm 3.1	30.4 \pm 8.1	55.7 \pm 4.8	18.3	15.7 \pm 0.18
		Bilinear	73.1 \pm 3.4	31.4 \pm 8.2	52.3 \pm 5.2	273.9	18.89 \pm 0.43
	ResDeepMixCNN	Concatenation	77.7 \pm 2.6	42.8 \pm 6.0	62.6 \pm 3.5	18.0	15.6 \pm 0.48
		Addition	80.3 \pm 2.9	48.0 \pm 7.3	64.4 \pm 4.3	11.5	15.65 \pm 0.31
		Attention-on-Mov	77.1 \pm 2.4	44.8 \pm 5.9	62.7 \pm 3.6	18.3	16.24 \pm 0.32
		Attention-on-Car	75.9 \pm 3.0	37.0 \pm 7.7	57.7 \pm 4.4	18.3	16.24 \pm 0.32
		Bilinear	72.8 \pm 3.2	29.8 \pm 6.5	52.1 \pm 4.4	273.9	18.88 \pm 0.4

Our observations corroborate a study using raw PSG signals for sleep stage classification [60]. That is using intermediate features instead of raw accelerometer data may alleviate the modality bias in the three-sleep stage classification task, while reducing the model parameters.

C THREE SLEEP STAGE CLASSIFICATION PERFORMANCE ON 21, 51 WINDOW LENGTH

C.1 The Effects of Sliding Windows Length

In addition to the window length of 101, we also conducted experiments based on the window lengths 51 and 21 followed the previous work [84]. For the Apple Watch dataset, the models with the highest mean F1, accuracy and Cohen's κ score in each fusion strategy were all based on the window length of 101. For the MESA dataset, we observed similar patterns on all feature settings. One possible explanation was when the time step of the input data became shorter, the intermediate features around the time point of the prediction might not contain enough information for three-stage sleep classification. This phenomenon corroborated the previous findings [84].

Table 8. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of the fusion strategy and method with the Apple Watch dataset using the ACT-HRS feature and evaluated at subject level during the recording period based on the window length of 51.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	73.1 \pm 3.1	40.5 \pm 6.7	58.7 \pm 4.4	9.1 \pm 22.5	-0.1 \pm 23.5	-9.0 \pm 8.8
	ResDeepCNN	Concatenation	75.4 \pm 2.9	43.9 \pm 6.4	61.1 \pm 3.9	23.0 \pm 24.1	-12.5 \pm 23.7	-10.5 \pm 7.2
Late-Stage Fusion	DeepCNN	Concatenation	74.1 \pm 2.7	41.9 \pm 6.1	59.8 \pm 3.7	12.0 \pm 21.6	-1.1 \pm 22.1	-10.9 \pm 7.9
		Addition	78.0 \pm 2.3	50.1 \pm 7.0	65.5 \pm 3.6	1.3 \pm 16.7	11.6 \pm 16.8	-12.9 \pm 6.6
Late-Stage Fusion	ResDeepCNN	Concatenation	76.6 \pm 2.5	45.9 \pm 6.9	62.6 \pm 4.2	20.1 \pm 18.8	-13.5 \pm 19.6	-6.7 \pm 7.8
		Addition	77.7 \pm 2.2	48.1 \pm 6.8	64.6 \pm 4.0	7.6 \pm 19.5	3.4 \pm 19.6	-11.0 \pm 5.8
Hybrid Fusion	DeepCNN	Concatenation	72.5 \pm 3.2	39.0 \pm 6.2	58.8 \pm 3.9	7.5 \pm 24.6	5.3 \pm 24.6	-12.8 \pm 6.8
		Addition	73.3 \pm 3.0	39.2 \pm 6.3	58.4 \pm 3.6	17.1 \pm 24.1	-0.8 \pm 24.1	-16.3 \pm 5.5
		Attention-on-Mov	72.6 \pm 3.0	39.0 \pm 6.0	59.4 \pm 3.4	15.1 \pm 23.6	-1.7 \pm 23.6	-13.4 \pm 6.6
		Attention-on-Car	72.7 \pm 2.9	37.2 \pm 6.4	58.3 \pm 3.7	23.7 \pm 21.0	-8.6 \pm 20.3	-15.2 \pm 6.1
		Bilinear	72.3 \pm 2.6	38.2 \pm 5.5	58.5 \pm 3.2	1.4 \pm 20.7	12.3 \pm 20.7	-13.7 \pm 6.3
	ResDeepCNN	Concatenation	73.6 \pm 2.9	41.5 \pm 6.5	60.7 \pm 4.1	8.8 \pm 23.1	1.4 \pm 24.9	-10.2 \pm 7.4
		Addition	73.1 \pm 3.1	40.5 \pm 6.5	59.7 \pm 3.9	11.5 \pm 23.8	-0.8 \pm 24.4	-10.7 \pm 6.8
		Attention-on-Mov	74.4 \pm 3.3	43.8 \pm 6.1	61.7 \pm 3.8	11.6 \pm 20.9	-1.5 \pm 21.5	-10.1 \pm 6.9
		Attention-on-Car	73.1 \pm 3.0	40.6 \pm 7.1	60.4 \pm 4.1	2.1 \pm 24.0	6.9 \pm 24.3	-9.0 \pm 7.9
		Bilinear	74.9 \pm 2.6	42.2 \pm 6.5	60.2 \pm 3.9	22.2 \pm 21.9	-12.2 \pm 22.0	-9.9 \pm 7.1

 Table 9. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategy and method with the Apple Watch dataset using the ACT-HRS feature and evaluated at subject level during the recording period based on the window length of 21.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	71.9 \pm 2.1	35.2 \pm 5.5	54.9 \pm 3.5	22.9 \pm 16.6	-13.9 \pm 16.8	-9.0 \pm 8.3
	ResDeepCNN	Concatenation	73.4 \pm 2.4	38.7 \pm 5.8	58.0 \pm 3.6	17.3 \pm 15.2	-9.6 \pm 16.0	-7.7 \pm 9.2
Late-Stage Fusion	DeepCNN	Concatenation	72.4 \pm 2.5	38.1 \pm 6.1	55.9 \pm 3.9	21.9 \pm 18.2	-7.4 \pm 18.0	-14.6 \pm 8.6
		Addition	75.3 \pm 2.4	39.7 \pm 7.7	59.2 \pm 4.2	14.6 \pm 13.9	-1.4 \pm 15.4	-13.1 \pm 7.2
Late-Stage Fusion	ResDeepCNN	Concatenation	72.7 \pm 2.6	38.3 \pm 6.8	57.3 \pm 4.2	11.9 \pm 20.3	-4.5 \pm 21.0	-7.4 \pm 8.8
		Addition	74.3 \pm 2.7	38.9 \pm 7.5	58.9 \pm 4.2	6.2 \pm 16.6	4.5 \pm 17.6	-10.7 \pm 7.9
Hybrid Fusion	DeepCNN	Concatenation	71.6 \pm 2.5	35.8 \pm 6.2	55.6 \pm 4.0	17.7 \pm 23.5	-0.0 \pm 23.4	-17.7 \pm 6.4
		Addition	71.6 \pm 2.6	35.0 \pm 5.7	55.8 \pm 3.6	22.9 \pm 19.7	-5.9 \pm 19.7	-17.0 \pm 6.7
		Attention-on-Mov	72.2 \pm 2.6	37.3 \pm 5.4	56.6 \pm 3.4	20.5 \pm 19.6	-3.6 \pm 20.8	-16.9 \pm 6.3
		Attention-on-Car	73.2 \pm 2.3	35.4 \pm 5.9	56.6 \pm 3.6	31.1 \pm 20.6	-13.3 \pm 20.1	-17.8 \pm 6.4
		Bilinear	71.5 \pm 2.9	37.3 \pm 6.0	57.2 \pm 3.7	5.0 \pm 19.1	6.9 \pm 17.7	-11.9 \pm 7.7
	ResDeepCNN	Concatenation	72.3 \pm 2.7	36.6 \pm 6.5	57.3 \pm 4.0	21.3 \pm 23.8	-5.3 \pm 23.3	-15.9 \pm 6.5
		Addition	71.6 \pm 2.2	35.2 \pm 5.2	55.5 \pm 3.4	24.7 \pm 18.2	-11.6 \pm 19.1	-13.2 \pm 6.9
		Attention-on-Mov	73.3 \pm 2.3	38.0 \pm 5.4	57.9 \pm 3.3	33.5 \pm 17.3	-19.3 \pm 17.7	-14.2 \pm 6.4
		Attention-on-Car	72.1 \pm 2.7	37.8 \pm 5.5	57.8 \pm 3.6	17.5 \pm 19.2	-5.8 \pm 18.8	-11.7 \pm 7.9
		Bilinear	70.5 \pm 2.7	35.7 \pm 5.8	55.9 \pm 3.6	3.1 \pm 19.2	1.2 \pm 19.5	-4.3 \pm 9.8

Table 10. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategy and method with the MESA test dataset using the ACT-HRS evaluated at subject level during the recording period based on the window length of 51.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	78.1 \pm 0.9	60.1 \pm 1.8	69.3 \pm 1.3	14.6 \pm 7.3	-18.7 \pm 3.6	4.1 \pm 6.9
	ResDeepCNN	Concatenation	76.7 \pm 1.0	60.2 \pm 1.9	70.3 \pm 1.3	-15.6 \pm 7.4	12.6 \pm 5.0	3.0 \pm 6.8
Late-Stage Fusion	DeepCNN	Concatenation	78.4 \pm 1.0	62.5 \pm 1.8	71.4 \pm 1.3	9.9 \pm 7.3	0.2 \pm 4.1	-10.0 \pm 6.4
		Addition	76.5 \pm 0.9	58.8 \pm 1.7	66.7 \pm 1.2	62.3 \pm 7.4	-21.7 \pm 3.7	-40.6 \pm 6.9
	ResDeepCNN	Concatenation	77.7 \pm 0.9	61.0 \pm 1.8	69.9 \pm 1.2	17.3 \pm 7.6	-4.6 \pm 4.2	-12.6 \pm 6.6
		Addition	74.8 \pm 1.0	55.7 \pm 1.8	66.2 \pm 1.3	-21.0 \pm 7.7	-16.0 \pm 4.2	37.0 \pm 7.5
Hybrid Fusion	DeepCNN	Concatenation	76.4 \pm 1.1	61.1 \pm 1.8	70.0 \pm 1.3	-9.4 \pm 8.0	17.9 \pm 5.0	-8.5 \pm 6.9
		Addition	77.2 \pm 1.0	61.3 \pm 1.7	70.6 \pm 1.2	-17.9 \pm 7.5	2.5 \pm 4.3	15.4 \pm 7.0
		Attention-on-Mov	74.6 \pm 1.1	59.1 \pm 1.8	69.0 \pm 1.2	-24.0 \pm 8.1	39.2 \pm 5.7	-15.3 \pm 6.6
		Attention-on-Car	77.8 \pm 0.9	60.7 \pm 1.8	69.8 \pm 1.2	5.4 \pm 7.4	-9.2 \pm 4.1	3.8 \pm 6.7
		Bilinear	77.1 \pm 0.9	60.1 \pm 1.8	70.2 \pm 1.2	6.8 \pm 8.0	13.9 \pm 5.0	-20.8 \pm 6.5
	ResDeepCNN	Concatenation	77.7 \pm 1.0	62.4 \pm 1.7	71.0 \pm 1.2	8.3 \pm 7.7	15.5 \pm 5.0	-23.8 \pm 6.3
		Addition	78.5 \pm 0.9	62.0 \pm 1.7	71.3 \pm 1.2	30.4 \pm 7.3	4.1 \pm 4.3	-34.5 \pm 6.5
		Attention-on-Mov	76.1 \pm 1.1	60.7 \pm 1.8	70.4 \pm 1.3	-9.9 \pm 8.2	31.3 \pm 5.7	-21.4 \pm 6.8
		Attention-on-Car	76.6 \pm 1.1	61.2 \pm 1.8	70.7 \pm 1.2	-0.3 \pm 7.6	25.8 \pm 5.1	-25.6 \pm 6.2
		Bilinear	76.7 \pm 0.9	60.2 \pm 1.7	69.7 \pm 1.2	-9.5 \pm 7.6	9.4 \pm 4.6	0.1 \pm 6.6

Table 11. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategy and method with the MESA test dataset using the ACT-HRS evaluated at subject level during the recording period based on the window length of 21.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	75.3 \pm 1.0	55.2 \pm 1.8	66.7 \pm 1.2	61.4 \pm 7.7	-4.3 \pm 4.2	-57.0 \pm 6.8
	ResDeepCNN	Concatenation	75.0 \pm 0.9	56.6 \pm 1.7	68.3 \pm 1.1	13.3 \pm 7.4	16.7 \pm 4.6	-30.0 \pm 6.5
Late-Stage Fusion	DeepCNN	Concatenation	76.6 \pm 0.9	57.8 \pm 1.6	68.0 \pm 1.2	29.8 \pm 7.0	-13.9 \pm 3.7	-15.9 \pm 6.4
		Addition	74.4 \pm 0.9	56.9 \pm 1.6	66.8 \pm 1.1	13.7 \pm 7.5	6.9 \pm 4.8	-20.7 \pm 6.4
	ResDeepCNN	Concatenation	75.9 \pm 0.9	58.1 \pm 1.6	68.8 \pm 1.1	7.8 \pm 7.1	6.8 \pm 4.1	-14.7 \pm 6.4
		Addition	74.7 \pm 0.9	56.6 \pm 1.6	66.6 \pm 1.1	33.4 \pm 7.4	2.7 \pm 4.7	-36.1 \pm 6.4
Hybrid Fusion	DeepCNN	Concatenation	76.0 \pm 0.9	57.1 \pm 1.7	68.6 \pm 1.2	33.8 \pm 7.6	8.2 \pm 4.4	-42.0 \pm 6.7
		Addition	74.8 \pm 1.0	55.7 \pm 1.8	67.9 \pm 1.2	22.6 \pm 7.8	21.7 \pm 4.9	-44.3 \pm 6.6
		Attention-on-Mov	74.2 \pm 1.0	56.8 \pm 1.7	67.9 \pm 1.2	-25.4 \pm 7.9	18.4 \pm 4.7	7.0 \pm 7.0
		Attention-on-Car	76.1 \pm 0.9	57.4 \pm 1.6	68.3 \pm 1.1	13.1 \pm 7.5	-1.0 \pm 4.3	-12.1 \pm 6.7
		Bilinear	76.9 \pm 0.9	57.6 \pm 1.7	68.0 \pm 1.2	26.4 \pm 7.2	-20.1 \pm 3.6	-6.3 \pm 6.8
	ResDeepCNN	Concatenation	76.5 \pm 0.9	57.6 \pm 1.7	68.6 \pm 1.2	47.2 \pm 7.4	-0.2 \pm 4.1	-47.0 \pm 6.6
		Addition	76.2 \pm 0.9	57.7 \pm 1.7	68.8 \pm 1.2	25.0 \pm 7.4	7.6 \pm 4.3	-32.6 \pm 6.6
		Attention-on-Mov	75.9 \pm 0.9	58.1 \pm 1.7	69.0 \pm 1.1	-0.2 \pm 7.5	10.0 \pm 4.3	-9.8 \pm 6.7
		Attention-on-Car	75.5 \pm 0.9	58.2 \pm 1.7	68.8 \pm 1.2	-5.1 \pm 7.5	13.2 \pm 4.6	-8.2 \pm 6.6
		Bilinear	77.0 \pm 0.9	58.7 \pm 1.6	68.6 \pm 1.1	39.9 \pm 7.1	-13.1 \pm 3.9	-26.8 \pm 6.5

Table 12. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategy and method in the MESA test dataset using the ACT-HRV feature set evaluated at subject level during the recording period based on the window length of 51.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	75.7 \pm 1.0	56.7 \pm 1.8	67.2 \pm 1.3	40.0 \pm 6.9	-11.1 \pm 3.9	-28.9 \pm 6.5
	ResDeepCNN	Concatenation	76.0 \pm 0.9	56.9 \pm 1.8	66.8 \pm 1.3	63.1 \pm 7.3	-17.3 \pm 3.6	-45.9 \pm 6.8
Late-Stage Fusion	DeepCNN	Concatenation	78.4 \pm 0.9	61.7 \pm 1.8	70.2 \pm 1.3	20.8 \pm 6.9	-10.9 \pm 3.7	-9.8 \pm 6.4
		Addition	77.6 \pm 0.9	60.5 \pm 1.8	69.2 \pm 1.2	34.3 \pm 6.9	-8.6 \pm 3.9	-25.7 \pm 6.5
Late-Stage Fusion	ResDeepCNN	Concatenation	78.0 \pm 1.0	62.4 \pm 1.7	71.1 \pm 1.2	2.8 \pm 7.2	3.3 \pm 4.1	-6.1 \pm 6.4
		Addition	77.5 \pm 0.9	61.1 \pm 1.8	70.3 \pm 1.2	16.8 \pm 6.9	2.4 \pm 4.0	-19.2 \pm 6.5
Hybrid Fusion	DeepCNN	Concatenation	77.2 \pm 1.0	60.3 \pm 1.7	69.7 \pm 1.3	12.5 \pm 7.5	-0.5 \pm 4.4	-12.0 \pm 6.5
		Addition	76.6 \pm 1.0	59.7 \pm 1.8	69.7 \pm 1.3	14.2 \pm 7.4	13.4 \pm 4.6	-27.6 \pm 6.5
		Attention-on-Mov	77.5 \pm 1.0	61.2 \pm 1.7	70.5 \pm 1.3	35.3 \pm 7.3	5.7 \pm 4.8	-40.9 \pm 6.4
		Attention-on-Car	77.8 \pm 0.9	60.4 \pm 1.7	68.5 \pm 1.3	45.2 \pm 7.1	-20.6 \pm 3.8	-24.7 \pm 6.4
		Bilinear	77.1 \pm 1.0	60.8 \pm 1.7	70.6 \pm 1.2	5.4 \pm 6.9	12.2 \pm 4.4	-17.7 \pm 6.2
	ResDeepCNN	Concatenation	76.7 \pm 1.1	60.7 \pm 1.9	70.6 \pm 1.3	6.2 \pm 7.9	23.4 \pm 5.4	-29.6 \pm 6.3
		Addition	76.7 \pm 1.1	61.3 \pm 1.8	70.9 \pm 1.3	-27.4 \pm 7.2	23.5 \pm 4.8	3.8 \pm 6.6
		Attention-on-Mov	78.7 \pm 0.9	62.2 \pm 1.7	70.0 \pm 1.3	43.9 \pm 6.9	-17.6 \pm 3.8	-26.3 \pm 6.4
		Attention-on-Car	78.3 \pm 0.9	62.2 \pm 1.7	70.6 \pm 1.3	37.7 \pm 7.1	-4.7 \pm 4.4	-33.0 \pm 6.2
		Bilinear	76.8 \pm 1.0	59.1 \pm 1.8	69.6 \pm 1.2	24.0 \pm 7.0	6.7 \pm 4.2	-30.6 \pm 6.2

 Table 13. Three-stage sleep classification results (mean \pm standard error at 95% confidence interval) for each combination of fusion strategy and method with the MESA test dataset using the ACT-HRV evaluated at subject level during the recording period based on the window length of 21.

Fusion Specifics			Performance Metrics			Time Deviation (min.)		
Fusion Strategy	Network	Fusion Method	Accuracy (%)	Cohen's κ	Mean F1 (%)	Non-REM sleep	REM sleep	Wake
Early-Stage Fusion	DeepCNN	Concatenation	75.9 \pm 0.9	57.0 \pm 1.7	67.5 \pm 1.2	34.7 \pm 7.0	-10.2 \pm 3.7	-24.6 \pm 6.6
	ResDeepCNN	Concatenation	75.4 \pm 0.9	56.0 \pm 1.7	67.2 \pm 1.2	8.5 \pm 7.1	-10.9 \pm 3.7	2.4 \pm 6.7
Late-Stage Fusion	DeepCNN	Concatenation	76.0 \pm 0.9	57.5 \pm 1.7	68.0 \pm 1.1	12.1 \pm 7.1	-4.5 \pm 3.9	-7.6 \pm 6.5
		Addition	74.4 \pm 0.9	54.4 \pm 1.7	64.0 \pm 1.2	26.7 \pm 7.7	-21.6 \pm 3.8	-5.1 \pm 7.2
Late-Stage Fusion	ResDeepCNN	Concatenation	76.2 \pm 0.9	58.2 \pm 1.7	68.0 \pm 1.2	16.7 \pm 6.9	-7.1 \pm 3.8	-9.6 \pm 6.5
		Addition	75.3 \pm 0.9	56.0 \pm 1.7	65.4 \pm 1.2	38.7 \pm 7.6	-19.2 \pm 3.8	-19.5 \pm 7.0
Hybrid Fusion	DeepCNN	Concatenation	75.8 \pm 1.0	57.5 \pm 1.7	68.5 \pm 1.2	21.2 \pm 7.7	3.6 \pm 4.4	-24.8 \pm 7.0
		Addition	75.2 \pm 1.0	57.0 \pm 1.8	68.5 \pm 1.2	11.3 \pm 7.4	21.4 \pm 5.2	-32.7 \pm 6.7
		Attention-on-Mov	75.4 \pm 1.0	56.7 \pm 1.7	67.7 \pm 1.2	15.3 \pm 7.5	3.9 \pm 4.9	-19.3 \pm 6.7
		Attention-on-Car	74.3 \pm 1.0	56.6 \pm 1.7	67.8 \pm 1.2	-14.6 \pm 7.5	14.7 \pm 4.6	-0.1 \pm 6.5
		Bilinear	74.9 \pm 1.0	57.4 \pm 1.8	68.3 \pm 1.2	-13.4 \pm 7.9	8.6 \pm 4.3	4.7 \pm 7.2
	ResDeepCNN	Concatenation	75.5 \pm 1.0	57.1 \pm 1.8	68.7 \pm 1.2	19.8 \pm 7.4	20.8 \pm 5.0	-40.6 \pm 6.6
		Addition	76.2 \pm 0.9	58.4 \pm 1.7	69.0 \pm 1.2	19.6 \pm 7.1	-0.9 \pm 4.1	-18.7 \pm 6.6
		Attention-on-Mov	76.0 \pm 1.0	58.7 \pm 1.8	69.0 \pm 1.2	18.5 \pm 7.4	10.3 \pm 4.9	-28.8 \pm 6.4
		Attention-on-Car	75.4 \pm 1.0	58.0 \pm 1.7	68.6 \pm 1.2	1.5 \pm 7.5	12.9 \pm 4.8	-14.4 \pm 6.5
		Bilinear	76.4 \pm 0.9	58.7 \pm 1.8	68.8 \pm 1.2	37.8 \pm 7.1	-0.8 \pm 4.1	-36.9 \pm 6.5