

PSG-MAE: Robust Multitask Sleep Event Monitoring using Multichannel PSG Reconstruction and Inter-channel Contrastive Learning

Yifei Wang, Qi Liu, *Senior Member, IEEE*, Fuli Min, and Honghao Wang

Abstract—Polysomnography (PSG) signals are essential for studying sleep processes and diagnosing sleep disorders. Analyzing PSG data through deep neural networks (DNNs) for automated sleep monitoring has become increasingly feasible. However, the limited availability of datasets for certain sleep events often leads to DNNs focusing on a single task with a single-sourced training dataset. As a result, these models struggle to transfer to new sleep events and lack robustness when applied to new datasets. To address these challenges, we propose PSG-MAE, a mask autoencoder (MAE) based pre-training framework. By performing self-supervised learning on a large volume of unlabeled PSG data, PSG-MAE develops a robust feature extraction network that can be broadly applied to various sleep event monitoring tasks. Unlike conventional MAEs, PSG-MAE generates complementary masks across PSG channels, integrates a multichannel signal reconstruction method, and employs a self-supervised inter-channel contrastive learning (ICCL) strategy. This approach enables the encoder to capture temporal features from each channel while simultaneously learning latent relationships between channels, thereby enhancing the utilization of multichannel information. Experimental results show that PSG-MAE effectively captures both temporal details and inter-channel information from PSG signals. When the encoder pre-trained through PSG-MAE is fine-tuned with downstream feature decomposition networks, it achieves an accuracy of 83.7% for sleep staging and 90.45% for detecting obstructive sleep apnea, which highlights the framework’s robustness and broad applicability.

Index Terms—Polysomnography Signal Analysis, Multichannel Signal Reconstruction, Pre-trained Deep Learning Models, Sleep Stage Classification, Obstructive Sleep Apnea Detection

I. INTRODUCTION

SLEEP is an essential necessity for life maintenance. Consistent and adequate rest is crucial for improving health, productivity, well-being and quality of life as well as public safety [1]. In recent years, the accelerated pace of global urbanization and rising stress have exacerbated the

This work was supported in part by the National Natural Science Foundation of China under Grant 62202174, in part by the Basic and Applied Basic Research Foundation of Guangzhou under Grant 2023A04J1674, in part by The Taihu Lake Innovation Fund for the School of Future Technology of South China University of Technology under Grant 2024B105611004, and in part by Guangdong Science and Technology Department Grant 2024A1313010012. (Corresponding author: Qi Liu and Honghao Wang.)

Yifei Wang and Qi Liu are with the School of Future Technology, South China University of Technology, Guangzhou 511400, China (e-mail: yswang634@outlook.com; drliuqi@scut.edu.cn)

Fuli Min and Honghao Wang are with the Department of Neurology, Guangzhou First People’s Hospital, and School of Medicine, South China University of Technology, Guangzhou 510180, China (e-mail: minfuli@163.com; wang_whh@163.com)

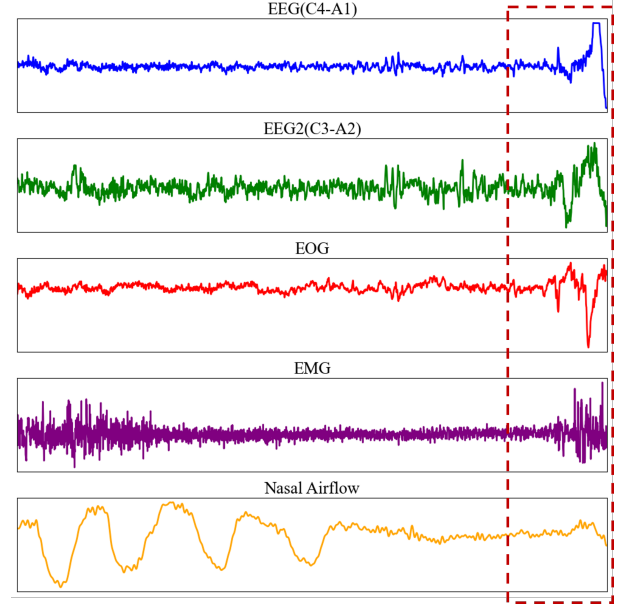


Fig. 1. Polysomnography (PSG) of one sleep epoch (30s), during which arousal occurs, marked in the dashed box. A sleep event during sleep often induces abrupt fluctuations in multiple channels of the PSG signals. Integrating the variations across different channels can help improve the accuracy of sleep event monitoring.

prevalence of sleep disorders, posing a substantial challenge to public health. Common sleep disorders, such as insomnia, arousal disorders, sleep apnea, rapid eye movement sleep behavior disorder (RBD), and periodic limb movement disorder (PLMD), are associated with a heightened risk of medical complications, including cardiovascular diseases, depression, and anxiety, diabetes and compromised immune function [2]. Consequently, the development of automated screening and intervention methods for sleep disorders is of significant research value.

Currently, combining polysomnography (PSG) with deep neural networks (DNNs) has become a widely explored approach in automated sleep monitoring research [3]. PSG is widely considered the most reliable method in the field of sleep medicine for diagnosing sleep-related disorders, often employed to evaluate both the diagnosis and efficacy of treatment for sleep disturbances [4], [5]. A standard PSG recording gathers data on brain waves (electroencephalography, EEG), eye movements (electrooculography, EOG), chin and leg muscle activity (electromyography, EMG), heart activity (elec-

trocardiography, ECG), chest and abdominal breathing effort, nasal airflow, oxygen saturation, etc [6]. While polysomnography (PSG) provides comprehensive documentation of sleep patterns, the analysis and clinical interpretation of these neurophysiological recordings demand rigorous systematic training. Additionally, annotating PSG data is labor-intensive and time-consuming, with 2-3 experts typically spending about 2 hours to annotate an 8-hour sleep recording. Subjective differences among experts can also lead to variability in annotation results [7], [8]. The process of annotating PSG sleep data includes the categorization of sleep stages and the identification of sleep events. In accordance with the sleep staging guidelines outlined by the American Academy of Sleep Medicine (AASM), PSG data is divided into 30-second segments along the temporal dimension. Each epoch is then classified into stages, including wake (W), non-rapid eye movement (NREM) stages (N1, N2, N3), and rapid eye movement (REM) sleep [9]. The epoch-based segmentation approach is utilized in contemporary clinical practice to label sleep events, including sleep apnea and limb movements, in order to ensure consistent analysis across research studies.

The multichannel nature of PSG signals makes them well-suited for integration with machine learning and deep neural networks, enabling the automatic extraction of complex sleep features and effectively modeling nonlinear relationships for accurate sleep event annotation [10], [11]. Current PSG-driven automated sleep events monitoring can be divided into two main areas. One focuses on the automatic sleep staging [12], [13], [14], while the other involves the detection and labeling of sleep behaviors, events, and disorders [15], [16]. However, there are two main challenges in current sleep event monitoring models. First, the wide variety of sleep events and their different manifestations across populations result in a limited quantity of public datasets for certain sleep events [17], [18]. As a result, many models are trained on small, task-specific datasets, making them sensitive to the feature distribution of the data. This limits their ability to transfer to other sleep event monitoring tasks and hinders a comprehensive, multidimensional evaluation of sleep. Second, most current sleep models rely on only a single or few PSG signal channels for specific monitoring tasks, neglecting the potential inter-channel interactions. As shown in Fig. 1, sleep events often induce signal changes across multiple channels. By integrating information from multiple channels, we can reduce misjudgments caused by disturbances in individual channels while improving the overall accuracy of event detection.

To address the challenges mentioned above, we propose leveraging unlabeled data through self-supervised learning to improve the model's stability in PSG feature extraction and enhance its performance in multichannel information fusion. To this end, we introduce PSG-MAE, a novel pre-training framework for PSG signals. In contrast to normal MAEs, PSG-MAE is based on a complementary-masking strategy for multichannel PSG signal reconstruction, meaning that one epoch (30 seconds) of PSG data is used as input with a pair of complementary masks generated along the channel dimension. Based on this design, we not only present a redesigned channel-level reconstruction loss but also introduce

inter-channel contrastive learning (ICCL) to further explore the inter-channel interaction information. PSG-MAE aims not only to capture fine-grained temporal information but also to learn the potential relationships between multiple channels of PSG. After the pre-training phase of PSG-MAE, a robust PSG encoder is built, which can be combined with downstream feature decomposition networks and fine-tuned to adapt to different sleep event monitoring tasks. The contributions of this paper can be summarized as follows:

- We propose PSG-MAE, a novel pre-training framework for PSG signals, which employs a complementary-masking strategy and leverages unlabeled PSG data for self-supervised learning. This approach enhances the feature extraction process, which is applicable to a wide range of sleep event monitoring tasks.
- To better exploit the multichannel nature of PSG signals, we introduce an updated channel-level signal reconstruction loss and a novel ICCL method. These innovations improve PSG-MAE's ability to capture fine-grained temporal information from each channel while also effectively modeling inter-channel interactions.
- The pre-trained PSG encoder demonstrates exceptional discriminative performance and robustness across multiple downstream sleep event monitoring tasks, including sleep staging and obstructive sleep apnea (OSA) detection, showcasing its broader applicability compared to traditional single-task models.

II. RELATED WORK

The research on PSG-driven automated sleep event monitoring currently focuses primarily on addressing the issues of sleep staging and the identification of other sleep events. Research on automated sleep staging generally follows two approaches: using raw signals as input and applying time-frequency transformations (e.g., Fourier transform, continuous wavelet transform) to generate spectrum as input [19]. Each channel of PSG signals is typically represented as one-dimensional time-series data. Consequently, feature extraction modules commonly employ models such as 1D convolutional neural networks (1D-CNNs) and recurrent neural networks (RNNs), which are well-suited for processing one-dimensional features and capturing temporal information. N. Goshtasbi et al. introduced SleepFCN, a fully convolutional framework that utilizes dual 1D-CNN branches with distinct kernel sizes to capture information across various EEG frequency bands [20]. Building upon SleepFCN, H. Zhu et al. developed MS-HNN, which integrates a squeeze-and-excitation (SE) block into the dual-kernel 1D-CNN branches to select more informative features. Additionally, MS-HNN employs bidirectional gated recurrent units (Bi-GRU) in the downstream network to learn temporal dependencies [21]. Y. Na et al. proposed using convolutional layers to fuse multichannel PSG data, allowing various physiological signals to contribute to the decision-making process [22]. Physiological signals, such as EEG, exhibit distinct variations in frequency bands and power distributions across different sleep stages. Time-frequency transformations effectively capture these frequency characteristics, especially

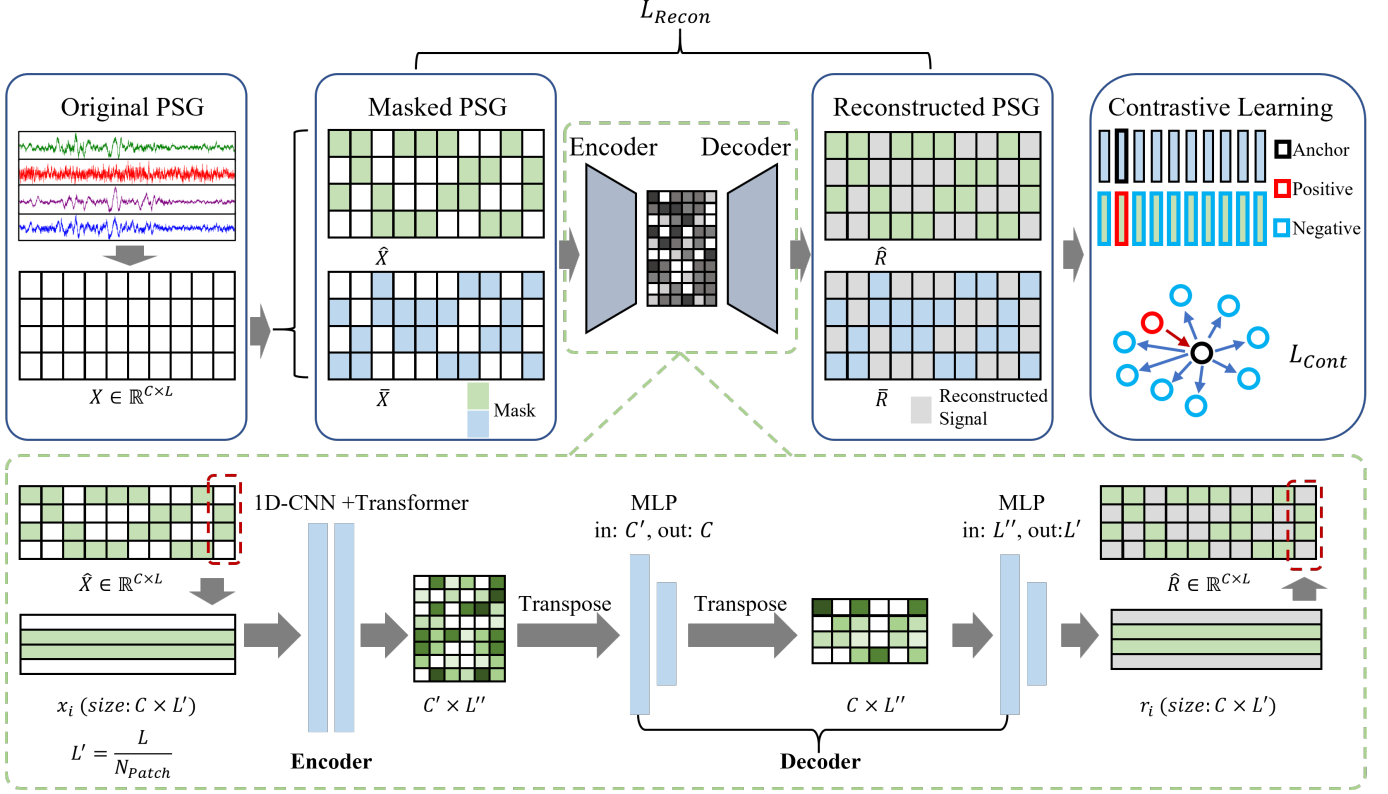


Fig. 2. The framework of **PSG-MAE**: The original PSG signal is divided into subsegments along the time dimension, followed by the application of complementary masks across the channel dimension. After passing through the encoder-decoder network, the unmasked portions of the signal are reconstructed, with the channel-level reconstruction loss facilitating the learning of temporal features in the original signal. In the pair of reconstructed PSG signals, one sub-segment is treated as an anchor, whose corresponding sub-segment in the other signal is considered as a positive sample, while the remaining subsegments are negative samples. ICCL is then applied to learn the intrinsic relationships between different channels by maximizing the distance of positive pairs and minimizing that of negative ones.

for applications involving non-stationary signals. P. Huy et al. introduced SeqSleepNet, which uses short-time Fourier transform (STFT) to convert PSG signals into time-frequency spectrograms and applies recurrent layers to capture both short-term and long-term dependencies within each epoch [23]. Y. Dai et al. proposed generating multichannel time-frequency spectrograms and employing multiple transformer groups to capture both individual channel features and joint features across channels [24].

Research on automated monitoring of other sleep events, such as sleep disorders, has advanced significantly. X. Zhao et al. segmented signals from the C3-A2 and C4-A1 EEG channels into five sub-bands, extracted entropy and variance features, and used machine learning to classify obstructive sleep apnea (OSA), central sleep apnea (CSA), and normal breathing events with [25]. A. Brink-Kjaer et al. used CNN+Bi-LSTM to extract features and temporal information from 5-minute PSG epochs for RBD classification, extending it with latent space transfer to analyze entire night recordings [26]. W. Qu et al. combined single-channel EEG data with a domain adaptation strategy, using similarity loss between encoders from source and target domains to learn temporal features. The source encoder was then integrated with LSTM networks for insomnia detection. [27].

Self-supervised learning has been shown to improve the

robustness of feature extraction by leveraging unlabeled data. In this context, MAE learns robust feature representations by masking and reconstructing portions of the input signals, and has demonstrated superior performance in various downstream tasks [28], [29], [30], [31]. MAE has been applied to the representation learning of temporal physiological signals. Y.-T. Lan et al. proposed the Corrupted Emotion Autoencoder (CEMOAE) framework to address channel corruption in EEG topographic maps by reconstructing masked signals to learn robust features and fine-tuning a pre-trained autoencoder for emotion recognition [32]. H. Ma et al. proposed a novel Region-State Masked Autoencoder (RS-MAE) that reduces redundancy in dynamic functional connectivity matrices, introduces region-state embeddings, and applies data augmentation to enhance classification performance for neuropsychiatric disorders based on resting-state fMRI. The encoder, pre-trained in this manner, has been shown to improve downstream task performance by capturing more relevant features [33].

III. METHODS

A. Overview

The general framework of **PSG-MAE**, based on the complementary masking strategy, is shown in Fig. 2. A 30-second multichannel PSG data segment, divided into temporally equal-length subsegments, serves as the input. After

that, the input is processed by applying a pair of randomly generated and complementary masks, forming two masked inputs that are fed into a shared encoder-decoder network to reconstruct the unmasked regions. The model then applies multichannel reconstruction and self-supervised ICCL to capture both temporal features and interactions among channels within the PSG input.

B. Multichannel Signal Reconstruction with Complementary-masking

The input PSG data $X \in \mathbb{R}^{C \times L}$ has C channels, with each channel containing L time steps. According to the standards of the International Classification of Sleep Disorders (ICSD), X encompasses a 30-second window of PSG data, where $L = 30s \times \text{sampling frequency}$. The input X is partitioned into smaller and manageable subsegments. These subsegments are defined by a hyperparameter N_{Patch} , which specifies how many patches the original time series data should be divided into along the time dimension. resulting in a set of subsegments x_i . This process can be expressed as

$$X = [x_1, x_2, x_3, \dots, x_N], \quad x_i \in \mathbb{R}^{C \times L'}, \quad L' = \frac{L}{N_{Patch}}. \quad (1)$$

At the beginning of the pre-training phase, the PSG-MAE framework generates a pair of complementary masks, M and $(1 - M) \in \mathbb{R}^{C \times L}$, each having the same size as the original input. The masking process begins by randomly selecting the floor of half the number of channels ($\lfloor C/2 \rfloor$) from each sub-segment x_i , and then combining all selected channels to form the mask M , while the channels that are not selected form the complementary mask $(1 - M)$. In this way, the two masks are complementary across the channel positions. Once the two complementary masks are created, the input matrix X is masked by applying both M and $(1 - M)$ across all subsegments, a pair of inputs to the shared encoder is generated as

$$\hat{X} = M * X, \quad (2)$$

$$\bar{X} = (1 - M) * X. \quad (3)$$

The masked data \hat{X} and \bar{X} , after undergoing the masking process, are passed through a transformer-based encoder for feature extraction and sequence modeling. This encoder leverages the self-attention mechanism to capture long-range dependencies within the data, enabling it to understand complex temporal patterns across multiple channels. After encoding, the transformed representation is fed into a multilayer perceptron (MLP)-based decoder. The decoder reconstructs the original signals from the encoded feature maps, producing a pair of reconstructed signals \hat{R} and \bar{R} . These reconstructed signals are subsequently compared with the original signals to compute the redesigned channel-level reconstruction loss, formulated as

$$L_{Recon} = L_{COS} + L_{MSE}, \quad (4)$$

$$L_{COS} = \frac{1}{C} \sum_{c=1}^C L_{channelCOS_c}, \quad (5)$$

$$L_{channelCOS_c} = 1 - \frac{1}{N} \sum_{n=1}^N \text{CosineSimilarity}(r_n^c, x_n^c), \quad (6)$$

$$\text{CosineSimilarity}(r_n^c, x_n^c) = \frac{\sum_{t=1}^{N_{patch}} r_n^c(t) x_n^c(t)}{\|r_n^c\| \|x_n^c\|}, \quad (7)$$

$$L_{MSE} = \frac{1}{C} \sum_{c=1}^C \frac{1}{T} \sum_{t=1}^T (x^c(t) - r^c(t))^2, \quad (8)$$

where reconstruction loss is formulated as a combination of channel-level cosine similarity loss L_{COS} and channel-level mean squared error (MSE) loss L_{MSE} . The L_{COS} is computed by firstly evaluating the cosine similarity of the c -th channel between the reconstructed signal subsegment r_n^c and the corresponding channel of the masked signal subsegment x_n^c , n means the n -th segment, as expressed in (7). The resulting cosine similarities are then averaged across both the subsegments and channels, as outlined in (6) and (5). The L_{MSE} is computed by first calculating the squared difference between the corresponding values of the masked signal $x^c(t)$ and the reconstructed signal $r^c(t)$ for each time step t along the channel c . The squared differences are then averaged over all T time steps and further averaged across all C channels, shown in (8). By integrating these two channel-level loss functions, the cosine similarity loss enforces the preservation of the overall pattern and trend of the reconstructed signals relative to the original data, while the MSE loss refines the relative magnitudes of the numerical values. This synergy between the two losses ensures that the model captures both the structural integrity and the numerical accuracy of the signals, leading to a faithful reconstruction of the original data.

C. Inter-channel Contrastive Learning

During sleep monitoring, the physiological data recorded by PSG is multi-source, encompassing signals collected from various sensors (such as EEG, EOG, respiratory airflow sensors, EMG, etc.) distributed across the body. These signal channels are typically synchronized in the time domain so that sleep events are often reflected simultaneously across multiple signal channels. This temporal alignment enables the multi-dimensional data to exhibit complementary information characteristics in multi-task sleep event detection. To fully exploit this synergy, PSG-MAE introduces a novel inter-channel contrastive learning (ICCL) strategy, which allows the encoder to uncover the latent commonalities and differences between signals, thereby enhancing the collaborative representation of features across different channels.

Specifically, we apply contrastive learning to the two groups of reconstructed output subsegments that form \hat{R} and \bar{R} . The objective is to ensure that signal blocks from different channels, which correspond to the same time frame, (i.e., originate from the same subsegment), are drawn closer together in the feature space. In contrast, channel blocks from different time frames, which have weak correlations, are pushed further apart. Upon obtaining the reconstructed signals \hat{R} and \bar{R} from the shared decoder, we recursively select subsegment \hat{r}_i from \hat{R} as anchor sample. Then the corresponding subsegment \bar{r}_i from

\bar{R} , which contains complementary channel information relative to \hat{r}_i within the same time interval, is designated as the positive sample. Meanwhile, the remaining subsegments $\hat{r}_{j \neq i}$ from \hat{R} serve as negative samples. During training, a triplet loss is employed to measure the relative distances among the anchor, positive, and negative samples. This strategy enables PSG-MAE to effectively learn and extract shared features across different channels while maintaining the independence and distinctiveness of temporal information. The channel contrastive loss L_{CL} is defined as

$$L_{CL} = \frac{1}{N_{patch}} \sum_{i=0}^{N_{patch}} F_{MAX}, \quad (9)$$

$$F_{MAX} = \max \left(0, d(\hat{r}_i, \bar{r}_i) - \frac{1}{N_{patch} - 1} \sum_{i \neq j} d(\hat{r}_i, \hat{r}_j) + \alpha \right), \quad (10)$$

$$d(x, y) = \sqrt{\sum_{k=1}^D (x_k - y_k)^2}, \quad (11)$$

where $d(x, y)$ denotes the Euclidean distance between samples x and y . x and y are two sample vectors with D dimensions. The components x_k and y_k represent the values of the samples in the k -th dimension. in (10), $d(x, y)$ measures the similarity between pairs of signal blocks, and α is a hyperparameter that defines the minimum margin between the positive and negative samples, ensuring that the negative samples are sufficiently far away from the anchor in the feature space.

The loss function for the pre-training framework for PSG signals based on a complementary-masking strategy for multichannel signal reconstruction is defined as

$$L = L_{Recon} + L_{CL}. \quad (12)$$

By jointly optimizing these two losses, PSG-MAE effectively preserves the fine-grained temporal details of PSG signals while also capturing the correlated features across different channels. The resulting encoded features provide a robust intermediate representation that can be leveraged for a wide array of downstream sleep-related tasks, including sleep stage classification, OSA detection, and other sleep events recognition. This enables the encoder to serve as a flexible and adaptable component that can be integrated into various temporal feature decomposition networks, each tailored to meet the specific requirements of sleep monitoring and classification. Consequently, the pre-trained encoder can be further fine-tuned for different sleep analysis applications, thereby enhancing model performance across diverse datasets and task-specific scenarios.

D. Downstream multitask sleep events monitoring

To effectively apply the pre-trained PSG-MAE encoder for downstream sleep event monitoring tasks, we design a feature decomposing network as illustrated in Fig. 3, this downstream network comprises several key components aimed at extracting and refining task-relevant information. To maximize the

utility of the pre-trained features derived from PSG-MAE, we incorporated a multi-branch 1D-CNN architecture. This architecture utilizes filters of varying sizes (1×3, 1×5, and 1×7) to capture multi-scale temporal patterns present in the PSG signals. These extracted features are then concatenated, enabling the network to integrate complementary information from different receptive fields. Subsequently, we apply a dimensionality-reduction step using a 1×1 filter, followed by global pooling, to further distill the feature representation while retaining the most informative components relevant to the task. The final output is then passed through an MLP layer for discrimination, generating task-specific results. A crucial aspect of this approach is the involvement of the pre-trained PSG-MAE encoder during the training phase of the downstream network. By integrating the encoder into the back-propagation process, the network can dynamically fine-tune its feature extraction capabilities. This enables the model to learn features that are specifically tailored to the downstream task. Rigorous validation of the effectiveness of the PSG-MAE encoder is conducted through following experiments on sleep staging and OSA detection tasks.

We employ cross-entropy loss to optimize the models for both downstream tasks. In the sleep staging task, multi-class cross-entropy loss (13) is used, while binary cross-entropy loss (14) is applied in the OSA detection task, they are defined as

$$L_{cls} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{ij} \log(p_{ij}), \quad (13)$$

$$L_{bcls} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)], \quad (14)$$

where y is the true label of sample, p is the sample's predicted probability of each class. Given the class imbalance in both sleep staging and OSA detection tasks (e.g., sleep apnea events typically constitute a small proportion of the total sleep duration), we apply class weights to the cross-entropy loss function, with higher weights assigned to underrepresented classes to mitigate the impact of class imbalance during training.

IV. EXPERIMENTS

The objective of this study is to achieve robust extraction of both single-channel temporal features and multichannel fusion features from PSG data by proposing the unsupervised learning approach, PSG-MAE. The experiments are designed with two main goals. First, to validate the effectiveness of the PSG-MAE, which uses complementary-masking and ICCL strategies to capture and reconstruct multichannel PSG signal information. Second, to assess the performance of the pre-trained encoder on two downstream sleep event monitoring tasks, sleep staging, and OSA detection, thereby evaluating its feasibility, applicability, and discriminative performance.

A. Dataset

To ensure sufficient diversity in the PSG data during the pre-training phase of the PSG-MAE and enhance the encoder's

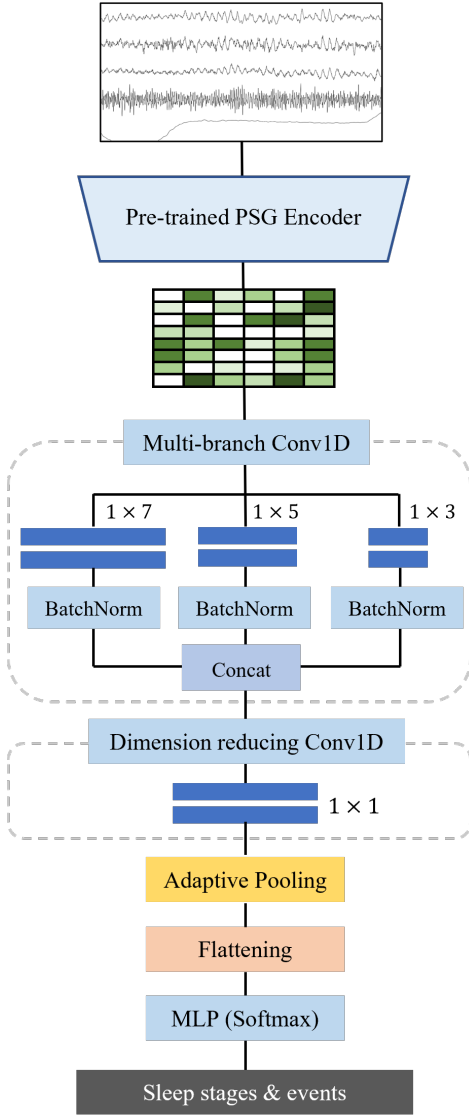


Fig. 3. Basic structure of downstream sleep events monitoring network.

robustness and generalization ability, this study uses three different datasets. These datasets not only help in improving the encoder’s performance but also allow for the evaluation of its effectiveness in downstream sleep event monitoring tasks:

The **Sleep Heart Health Study (SHHS)** [34], [35] is a multicenter epidemiological research resource aimed at assessing the impact of sleep-disordered breathing on cardiovascular health and other health outcomes. The dataset originates from a study led by the National Heart, Lung, and Blood Institute (NHLBI), with participants from various communities across the United States. It includes approximately 6,000 adults, primarily aged 50 and older, with data collection beginning in 1995 and ongoing long-term follow-up. The dataset covers multiple physiological signal channels with a sampling frequency of 100Hz.

The **PSG-audio** dataset [36] is sourced from the Sismanoglu – Amalia Fleming General Hospital in Athens, Greece, and was collected and annotated by the hospital’s medical team. The dataset contains 212 synchronized PSG recordings,

which also include audio recordings of breathing sounds from both tracheal and ambient microphones, for the analysis of apneic events and the development of home-based apnea detection techniques. The sampling frequency of the EEG signals is 200Hz.

The **Clinical-PSG** dataset was collected by the Department of Neurology at Guangzhou First People’s Hospital in Guangdong, China, within a clinical laboratory setting. This It contains PSG data recorded from 371 subjects during their nighttime sleep from 2021 to 2022. The EEG channels have a sampling frequency of 200 Hz, and sleep events are annotated in 30-second sleep epochs by professional sleep medicine physicians. These annotations include sleep stages, apneas, PLMD, and periodic limb movements while awake (PLMA). This dataset is intended for the diagnosis and research of sleep disorders. The construction and utilization processes of this dataset involved no collection of privacy information from subjects and were approved by the Guangzhou First People’s Hospital Ethics Committee (Approval No. K-2025-067-01).

B. Experimental Setup

During the pre-training phase, the synthetic dataset comprises 2,200 nights of PSG data from the SHHS, PSG-audio, and Clinical-PSG datasets. Combining these datasets ensures that the encoder learns the distribution of physiological signals across different datasets and subjects, enhancing the robustness of feature extraction. To maintain consistency in channel selection across the pre-training data, five common channels are chosen from the datasets: right central brain activity, left central brain activity, left eye movement, jaw electromyography, and pressure-based airflow signal. The channel names selected and the amount of data used in the pre-training phase from each dataset are shown in Table I. The PSG data undergoes EEG artifact removal processing to reduce interference from non-brain signals, thus enhancing the quality and reliability of the data. The sampling frequency of the PSG data is set to 100Hz, therefore the shape of the processed sleep epochs, (*channel number*, *data length*) is standardized to (5, 3000). To mitigate individual differences in the PSG data, we apply the Z-Score normalization method to each of the five channels in the PSG signals, the formula is

$$x' = \frac{x - \mu}{\sigma}. \quad (15)$$

This approach involves calculating the median (μ) and standard deviation (σ) for each channel, and then using these values to normalize the data accordingly. The training data is randomly shuffled, with 80% allocated for training, 10% used for validation, and the remaining 10% used for testing.

In subsequent downstream task experiments, the SHHS dataset is utilized for the sleep staging task, where 600 nights of data, not included in the pre-training phase, are randomly selected. In contrast, the remaining 171 nights from the PSG-clinical dataset, with detailed annotations, are employed for OSA detection. The data used in both tasks retains the same channel selection as in the pre-training phase to maintain consistency across tasks. Given the considerable subject-level

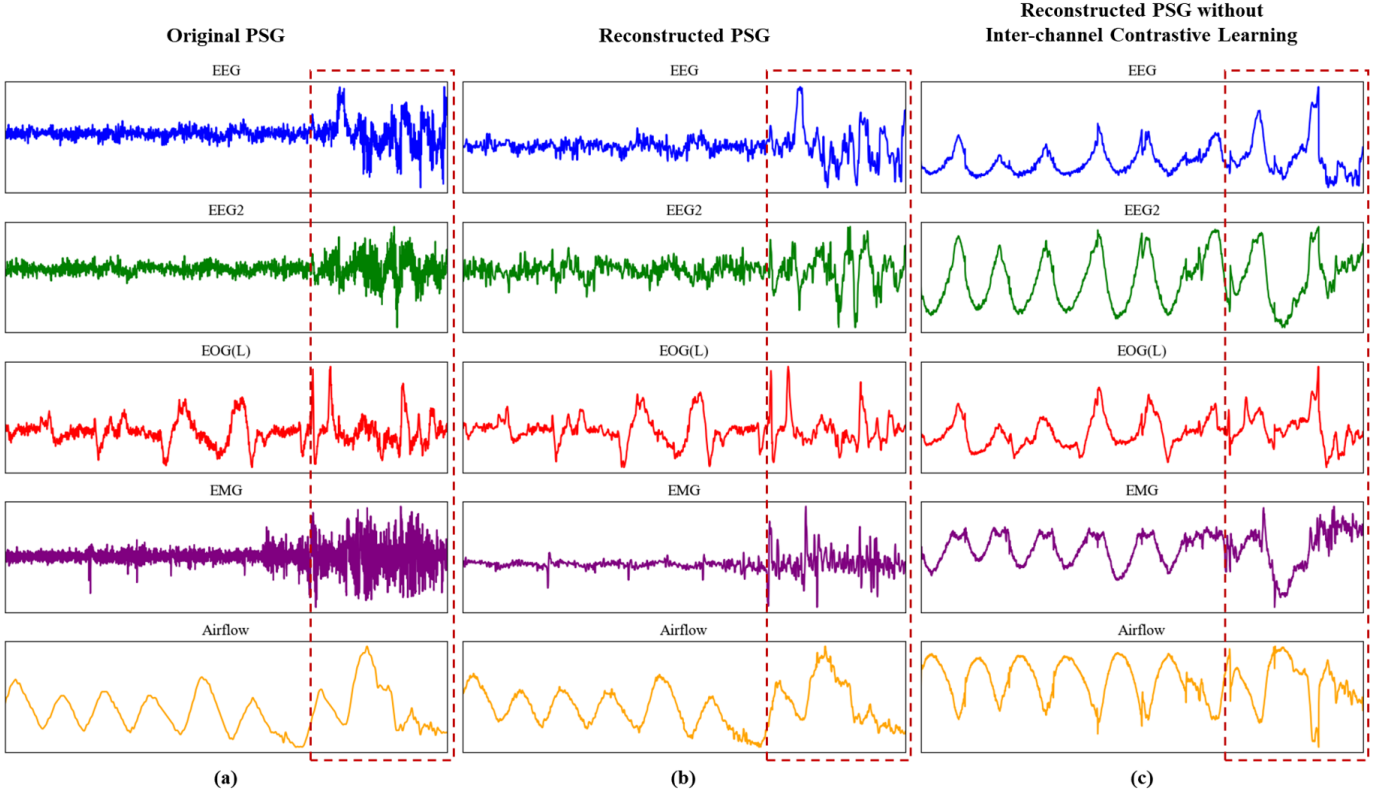


Fig. 4. The signal reconstruction results of PSG-MAE pre-training, show that the channel-level signal reconstruction loss and ICCL enable the framework to learn fine-grained temporal information within PSG channels as well as interaction information between channels. In contrast, without ICCL, it becomes difficult to disentangle individual channel information from the fused multichannel features.

TABLE I
CHANNEL CHOSEN FROM PSG DATASETS

Datasets	Data Number	PSG Channels
SHHS	1900	'EEG (C4-A1)', 'EEG2 (C3-A2)', 'EOG (L)', 'EMG' and 'AIRFLOW'
PSG-audio	200	'EEG C3-A2', 'EEG C4-A1', 'EOG LOC-A2', 'EMG Chin' and 'Flow Patient (Pressure cannula)'
PSG-clinical	100	'EEG C3-REF', 'EEG C4-REF', 'EOG LOC', 'EMG Chin' and 'Airflow'

variability inherent in PSG signals, both downstream tasks employ a subject-wise 5-fold cross-validation strategy. Specifically, the PSG data are randomly partitioned at the subject level into five equal subsets. In each fold, 80% of the data from four subsets is used for training, 20% for validation, and the remaining subset is reserved for testing. The final experimental results are obtained by averaging the outcomes across the five folds.

C. Evaluation Metrics

To rigorously evaluate the effectiveness of the upstream pre-training phase of PSG-MAE and its performance in downstream sleep events monitoring tasks, this study firstly employs the mean squared error (MSE) to quantify the discrepancy between the reconstructed and original signals, thereby assessing the accuracy of the signal reconstruction. The MSE is written as

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_{\text{pred}} - y_{\text{true}})^2. \quad (16)$$

Subsequently, in the sleep staging task, once the true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN) are computed for each class, accuracy (ACC), precision, recall and macro F1-score (MF1) are calculated as

$$ACC = \frac{TP + TN}{TP + TN + FP + FN}, \quad (17)$$

$$Precision = \frac{TP}{TP + FP}, \quad (18)$$

$$Recall = \frac{TP}{TP + FN}, \quad (19)$$

$$F1\text{-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}. \quad (20)$$

Due to the infrequent and brief nature of OSA events, the positive class is underrepresented in the OSA detection

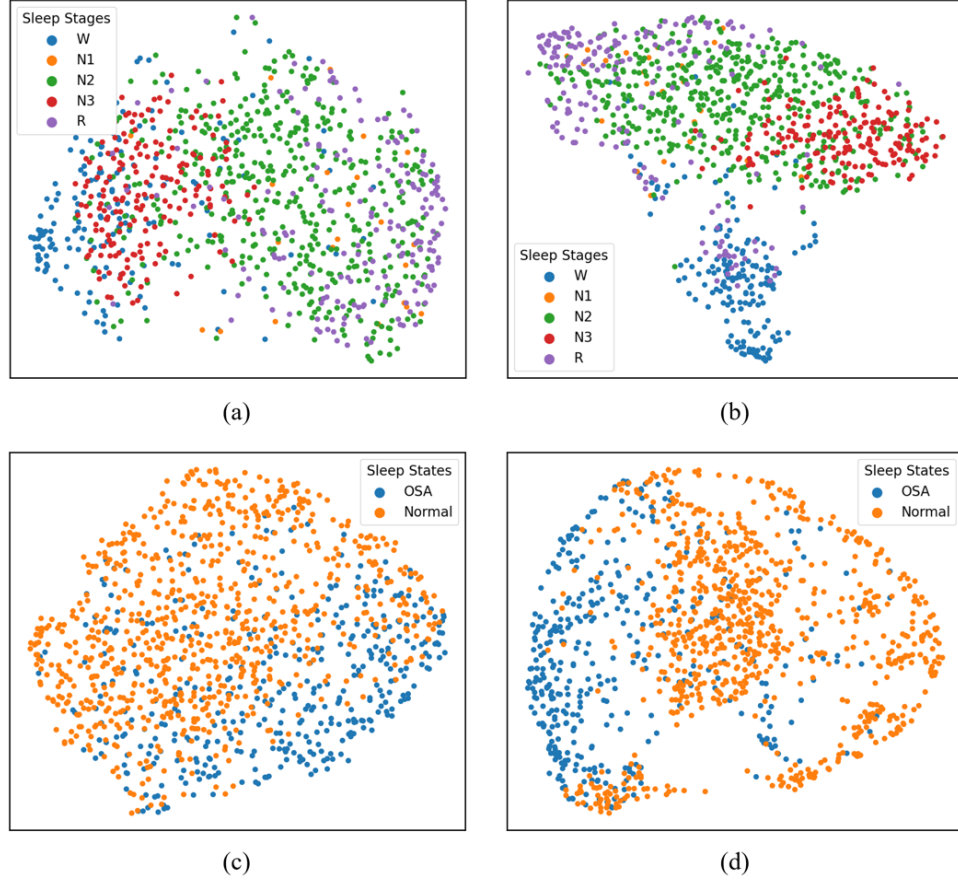


Fig. 5. UMAP visualization of PSG-MAE extracted features before and after downstream task training: (a) and (b) for sleep staging, (c) and (d) for OSA detection.

task. As a result, accuracy alone may not provide a reliable evaluation. Therefore, MF1 is used alongside accuracy, as this metric offers a more balanced assessment of the model's performance and helps mitigate the bias introduced by the dominant negative samples.

V. EXPERIMENTAL RESULTS & DISCUSSION

A. Results of PSG-MAE Pre-training Process

The Fig 4. shows the original PSG signal input and the reconstructed PSG signal output from PSG-MAE during the validation of pre-training phase. A comparison between (a) and (b) shows that the reconstructed signal accurately replicates the trend of signal variations while preserving the temporal details of the original signal. Furthermore, in the highlighted dashed boxes, fluctuations of multiple signal channels in the original signal are aligned with the time scale of the reconstructed signal, indicating that PSG-MAE is capable of capturing underlying relationships among different channels. Through ICCL, PSG-MAE retains the original signal's trend while demonstrating a certain degree of noise suppression, which indicates the interaction content learned by the encoder can help repair corrupted signal channels. The comparison between (2) and (3) underscores the critical role of ICCL in preserving single-channel information within multi-channel fused features. Without ICCL, PSG-MAE struggles to accurately reconstruct signals due to the loss of channel-specific

details. This highlights ICCL's ability to enhance the fusion of multi-channel information during the pre-training phase while minimizing information loss from individual channels. Table II presents the MSE of the selected 5 channels from the reconstructed and original signals, quantitatively reflecting the accuracy of PSG-MAE in signal reconstruction. EEG and EOG signal reconstruction demonstrate good performance, with effective preservation of waveform details and favorable MSE values. Similarly, the PSG-MAE model without ICCL cannot achieve the same level of reconstruction performance. Despite the high noise levels in the original EMG signals, the PSG-MAE model is still able to recover the main trends of the signal. In contrast, airflow signals show greater variability due to the significant time gap between the datasets and variations in data acquisition conditions, which complicate the reconstruction process. The reconstruction performance of some airflow signals is suboptimal, resulting in relatively higher average MSE values.

B. Downstream Sleep Event Monitoring Results

The Fig 5. shows the uniform manifold approximation and projection (UMAP) visualization of features extracted by the pre-trained encoder. Each point represents the reduced feature of one sleep epoch. The feature distribution before and after the training of the sleep staging task is shown in (a) and

TABLE II
THE MSE VALUE OF PSG CHANNELS BETWEEN THE ORIGINAL SIGNAL AND THE RECONSTRUCTED SIGNAL OF PSG-MAE.

PSG Channels	EEG	EEG2	EOG (L)	EMG	Airflow
MSE without ICCL	2.7×10^{-2}	2×10^{-2}	8.9×10^{-2}	1.5×10^{-2}	1.69
MSE with ICCL	8×10^{-6}	7×10^{-6}	3.6×10^{-5}	6×10^{-6}	7.66×10^{-2}

TABLE III
METHODS AND PERFORMANCE METRICS OF SLEEP STAGING TASK

Methods	Accuracy	Macro F1-score	ACC. W	ACC. N1	ACC. N2	ACC. N3	ACC. R
SleepEEGNet [37]	73.9%	68.4%	81.3%	34.4%	73.4%	75.9%	77.0%
DeepSleepNet [38]	81.0%	73.9%	85.4%	40.5%	82.5%	79.3%	81.9%
MultitaskCNN [39]	81.4%	71.2%	82.2%	25.7%	85.5%	83.3%	81.1%
AttnSleep [40]	82.3%	74.1%	85.0%	34.2%	85.7%	83.5%	82.3%
CausalAttenNet [41]	83.3%	73.1%	85.4%	27.6%	84.8%	84.0%	82.9%
Ours	83.7%	74.7%	94.8%	33.2%	84.5%	85.3%	79.8%

(b) respectively. A comparison reveals that, before the downstream task training, the features extracted by the pre-trained encoder exhibit only subtle clustering. After the downstream task training, features corresponding to different sleep stages exhibit more distinct clustering, indicating that the encoder has been fine-tuned to better align with the feature extraction requirements of the downstream task. Table III provides the sleep staging performance of the fine-tuned encoder, coupled with the downstream classification head. Compared to other approaches, the proposed PSG-MAE framework demonstrates superior performance in terms of prediction accuracy and MF1 score, showing enhanced discrimination, particularly for the W (waking) and N3 (deep sleep) stages.

TABLE IV
METHODS AND PERFORMANCE METRICS FOR OSA DETECTION

Methods	Accuracy	Macro F1-score
RF	83.7%	56.1%
SVM	87.1%	46.5%
CNN	87.6%	46.7%
Ours	90.45%	67.33%

The UMAP visualization in (c) and (d) of OSA detection training on the PSG-clinical dataset shows that features of OSA epochs and normal epochs are intermixed in the feature space before training and then a more distinct separation between the two classes emerges after the downstream fine-tuning of the encoder, indicating that the fine-tuned encoder is capable of distinguishing between OSA and normal sleep states. Table IV presents the performance of the fine-tuned encoder combined with the downstream network for OSA detection. The model is compared to the machine learning

methods of random forest (RF) and support vector machine (SVM) and a normal 1D-CNN network. These comparative models also take the raw PSG signals as input and the proposed approach shows a clear advantage in diagnostic accuracy and the MF1 score reaches its optimal value, demonstrating that the proposed network can effectively detect OSA even in the case of imbalanced label distribution.

C. Discussion

Experimental results validate the effectiveness of the PSG-MAE pre-training process and its feasibility for adaptation to downstream sleep event monitoring tasks. Specifically, by adding a pair of complementary masks to multiple channels of PSG segments from the same time subsegment, the channel-level signal reconstruction loss ensures that the encoder can learn to extract temporal features from multiple channels in the signal reconstruction process. Additionally, the complementary channels engage in ICCL with channels from different time segments, enabling the encoder to capture inter-channel interaction information at the same time point. Experimental results demonstrate that, for noisy channels, ICCL helps clarify the trend of reconstructed signal variations, and ensure simultaneously reconstructed signals involving multiple channels can preserve channel-related information and align with time steps. During the training of downstream sleep event monitoring tasks, the pre-trained encoder, when combined with the downstream feature decomposing network and fine-tuned, shows strong performance in both sleep staging and OSA detection. This indicates that PSG-MAE is not only suitable for the current tasks but also has the potential to adapt to a variety of downstream tasks, thereby facilitating the development of a multi-dimensional system for comprehensive sleep assessment.

VI. CONCLUSION

We propose an MAE-based pre-training framework named PSG-MAE, which leverages unlabeled PSG data through self-supervised learning to enhance the feature extraction capability of automated sleep event monitoring networks. In the pre-training phase, the channel-level signal reconstruction loss ensures the extraction of fine-grained time-series features, while ICCL emphasizes channel interaction information. The training process adopts PSG data from different datasets with the same channel configuration to improve adaptability to diverse data distributions. By integrating with downstream task networks, the pre-trained encoder can be fine-tuned to perform multitask sleep event monitoring. Compared to traditional single-task sleep monitoring networks, PSG-MAE demonstrates greater versatility. Experimental results show that the PSG-MAE pre-training framework can learn stable PSG features and achieve remarkable performance in both sleep staging and OSA detection tasks. The current limitation of PSG-MAE framework lies in the suboptimal performance of fine-tuned sleep staging network in recognizing N1 and R sleep stages. This issue may stem from the relatively small proportion of N1 and R stages in the training dataset and the lack of targeted strategies for data augmentation. In the future, efforts will be directed toward refining the downstream network and training process to improve sleep staging accuracy. Additionally, the framework will be expanded to support a broader range of automated sleep event monitoring tasks. The ultimate objective is to develop a multidimensional sleep monitoring system capable of delivering a more comprehensive analysis of PSG signals from a single input segment.

REFERENCES

REFERENCES

- [1] K. Ramar, R. K. Malhotra, K. A. Carden, J. L. Martin, F. Abbasi-Feinberg, R. N. Aurora, V. K. Kapur, E. J. Olson, C. L. Rosen, J. A. Rowley *et al.*, "Sleep is essential to health: an american academy of sleep medicine position statement," *Journal of Clinical Sleep Medicine*, vol. 17, no. 10, pp. 2115–2119, 2021.
- [2] M. Lee, H. Kang, S.-H. Yu, H. Cho, J. Oh, G. Van Der Lande, O. Gossieres, and J.-H. Jeong, "Automatic sleep stage classification using nasal pressure decoding based on a multi-kernel convolutional bilstm network," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 2024.
- [3] M. Yazdi, M. Samaee, and D. Massicotte, "A review on automated sleep study," *Annals of Biomedical Engineering*, vol. 52, no. 6, pp. 1463–1491, 2024.
- [4] J. V. Rundo and R. Downey, "Chapter 25 - polysomnography," in *Clinical Neurophysiology: Basis and Technical Aspects*, ser. Handbook of Clinical Neurology, K. H. Levin and P. Chauvel, Eds. Elsevier, 2019, vol. 160, pp. 381–392. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780444640321000254>
- [5] H. Zhou, A. Liu, S. Ding, J. Yao, and X. Chen, "An interpretable single-channel eeg sleep staging model based on prototype matching and multi-task learning," *IEEE Sensors Journal*, 2024.
- [6] H. Zhang, X. Wang, H. Li, S. Mehendale, and Y. Guan, "Auto-annotating sleep stages based on polysomnographic data," *Patterns*, vol. 3, no. 1, 2022.
- [7] I. Perez-Pozuelo, B. Zhai, J. Palotti, R. Mall, M. Aupetit, J. M. Garcia-Gomez, S. Taheri, Y. Guan, and L. Fernandez-Luque, "The future of sleep health: a data-driven revolution in sleep science and medicine," *NPJ digital medicine*, vol. 3, no. 1, p. 42, 2020.
- [8] E. Khalili and B. M. Asl, "Automatic sleep stage classification using temporal convolutional neural network and new data augmentation technique from raw single-channel eeg," *Computer Methods and Programs in Biomedicine*, vol. 204, p. 106063, 2021.
- [9] R. Berry, R. Brooks, C. Gamaldo, S. Harding, R. Lloyd, C. Marcus, B. Vaughn, and A. A. of Sleep Medicine, *The AASM Manual for the Scoring of Sleep and Associated Events: Rules, Terminology and Technical Specifications : Version 2.3*. American Academy of Sleep Medicine, 2015. [Online]. Available: <https://books.google.com/books?id=SySXAQAACAAJ>
- [10] J. Somanna, D. Joshi, H. Gundu, and G. Srinivasa, "Automated classification of sleep apnea and hypopnea on polysomnography data," in *2019 12th Biomedical Engineering International Conference (BMEiCON)*. IEEE, 2019, pp. 1–5.
- [11] A. De and E. Priya, "Sleep apnea sub-type detection from polysomnography signals," in *2024 IEEE International Conference on Interdisciplinary Approaches in Technology and Management for Social Innovation (IATMSI)*, vol. 2. IEEE, 2024, pp. 1–6.
- [12] S. K. Satapathy, S. Thakkar, A. Patel, and D. Patel, "A machine learning-based models for intelligent automated sleep staging classification system using polysomnography data," in *2023 IEEE 11th Region 10 Humanitarian Technology Conference (R10-HTC)*. IEEE, 2023, pp. 267–272.
- [13] A. Procházka, J. Kuchyňka, M. Yadollahi, C. P. S. Araujo, and O. Vyšata, "Adaptive segmentation of multimodal polysomnography data for sleep stages detection," in *2017 22nd International Conference on Digital Signal Processing (DSP)*. IEEE, 2017, pp. 1–4.
- [14] D. Zhang, Y. She, J. Sun, X. Yang, X. Zeng, and W. Qin, "Swinsleep: A deep learning framework advancing overnight sleep staging towards clinical practice," *IEEE Transactions on Instrumentation and Measurement*, 2024.
- [15] X. Li, A. Al-Ani, and S. H. Ling, "Feature selection for the detection of sleep apnea using multi-bio signals from overnight polysomnography," in *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE, 2018, pp. 1444–1447.
- [16] A. Bartolo, B. D. Clymer, R. C. Burgess, J. P. Turnbull, J. A. Golish, and M. C. Perry, "An arrhythmia detector and heart rate estimator for overnight polysomnography studies," *IEEE transactions on biomedical engineering*, vol. 48, no. 5, pp. 513–521, 2001.
- [17] F. Ehrlich, T. Sehr, M. Brandt, M. Schmidt, H. Malberg, M. Sedlmayr, and M. Goldammer, "State-of-the-art sleep arousal detection evaluated on a comprehensive clinical dataset," *Scientific Reports*, vol. 14, no. 1, p. 16239, 2024.
- [18] H. Lee, B. Li, S. DeForte, M. L. Splaingard, Y. Huang, Y. Chi, and S. L. Linwood, "A large collection of real-world pediatric sleep studies," *Scientific Data*, vol. 9, no. 1, p. 421, 2022.
- [19] R. N. Sekkal, F. Bereksi-Reguig, D. Ruiz-Fernandez, N. Dib, and S. Sekkal, "Automatic sleep stage classification: From classical machine learning methods to deep learning," *Biomedical Signal Processing and Control*, vol. 77, p. 103751, 2022.
- [20] N. Goshtasbi, R. Boostani, and S. Saneii, "Sleepfcn: A fully convolutional deep learning framework for sleep stage classification using single-channel electroencephalograms," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 30, pp. 2088–2096, 2022.
- [21] H. Zhu, L. Wang, N. Shen, Y. Wu, S. Feng, Y. Xu, C. Chen, and W. Chen, "Ms-hnn: Multi-scale hierarchical neural network with squeeze and excitation block for neonatal sleep staging using a single-channel eeg," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 31, pp. 2195–2204, 2023.
- [22] Y. Na, D. Kim, D.-K. Kim, and J.-G. Lee, "Evaluation of osa patient sleep stage classification performance using a multi-channel psg dataset," in *2022 IEEE International Conference on Consumer Electronics-Asia (ICCE-Asia)*. IEEE, 2022, pp. 1–4.
- [23] P. Huy, F. Andreotti, N. Cooray, O. Y. Chen, and M. De Vos, "Seqsleepnet: end-to-end hierarchical recurrent neural network for sequence-to-sequence automatic sleep staging," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 27, no. 3, pp. 400–410, 2019.
- [24] Y. Dai, X. Li, S. Liang, L. Wang, Q. Duan, H. Yang, C. Zhang, X. Chen, L. Li, X. Li *et al.*, "Multichannelsleepnet: A transformer-based model for automatic sleep stage classification with psg," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 9, pp. 4204–4215, 2023.
- [25] X. Zhao, X. Wang, T. Yang, S. Ji, H. Wang, J. Wang, Y. Wang, and Q. Wu, "Classification of sleep apnea based on eeg sub-band signal characteristics," *Scientific Reports*, vol. 11, no. 1, p. 5824, 2021.
- [26] A. Brink-Kjaer, K. M. Gunter, E. Mignot, E. Doring, P. Jennum, and H. B. Sorensen, "End-to-end deep learning of polysomnograms for classification of rem sleep behavior disorder," in *2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 2022, pp. 2941–2944.

- [27] W. Qu, C.-H. Kao, H. Hong, Z. Chi, R. Grunstein, C. Gordon, and Z. Wang, "Single-channel eeg based insomnia detection with domain adaptation," *Computers in biology and medicine*, vol. 139, p. 104989, 2021.
- [28] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 16 000–16 009.
- [29] L. Wang, B. Huang, Z. Zhao, Z. Tong, Y. He, Y. Wang, Y. Wang, and Y. Qiao, "Videomae v2: Scaling video masked autoencoders with dual masking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 14 549–14 560.
- [30] Z. Cai, S. Ghosh, K. Stefanov, A. Dhall, J. Cai, H. Rezatofighi, R. Haffari, and M. Hayat, "Marlin: Masked autoencoder for facial video representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 1493–1504.
- [31] Z. Zhang, P. Zhao, E. Park, and J. Yang, "Mart: Masked affective representation learning via masked temporal distribution distillation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 12 830–12 840.
- [32] Y.-T. Lan, W.-B. Jiang, W.-L. Zheng, and B.-L. Lu, "Cemoae: A dynamic autoencoder with masked channel modeling for robust eeg-based emotion recognition," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 1871–1875.
- [33] H. Ma, Y. Xu, and L. Tian, "Rs-mae: Region-state masked autoencoder for neuropsychiatric disorder classifications based on resting-state fmri," *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [34] G.-Q. Zhang, L. Cui, R. Mueller, S. Tao, M. Kim, M. Rueschman, S. Mariani, D. Mobley, and S. Redline, "The national sleep research resource: towards a sleep data commons," *Journal of the American Medical Informatics Association*, vol. 25, no. 10, pp. 1351–1358, 2018.
- [35] S. F. Quan, B. V. Howard, C. Iber, J. P. Kiley, F. J. Nieto, G. T. O'Connor, D. M. Rapoport, S. Redline, J. Robbins, J. M. Samet, and P. W. Wahl, "The sleep heart health study: Design, rationale, and methods," *Sleep*, vol. 20, no. 12, pp. 1077–1085, 12 1997. [Online]. Available: <https://doi.org/10.1093/sleep/20.12.1077>
- [36] G. Korompili, A. Amfilochiou, L. Kokkalas, S. A. Mitilneos, N.-A. Tatlas, M. Kouvaras, E. Kastanakis, C. Maniou, and S. M. Potirakis, "Psg-audio, a scored polysomnography dataset with simultaneous audio recordings for sleep apnea studies," *Scientific data*, vol. 8, no. 1, p. 197, 2021.
- [37] S. Mousavi, F. Afghah, and U. R. Acharya, "Sleeppegnet: Automated sleep stage scoring with sequence to sequence deep learning approach," *PloS one*, vol. 14, no. 5, p. e0216456, 2019.
- [38] A. Supratak, H. Dong, C. Wu, and Y. Guo, "Deepsleepnet: A model for automatic sleep stage scoring based on raw single-channel eeg," *IEEE transactions on neural systems and rehabilitation engineering*, vol. 25, no. 11, pp. 1998–2008, 2017.
- [39] H. Phan, F. Andreotti, N. Cooray, O. Y. Chén, and M. De Vos, "Joint classification and prediction cnn framework for automatic sleep stage classification," *IEEE Transactions on Biomedical Engineering*, vol. 66, no. 5, pp. 1285–1296, 2019.
- [40] E. Eldele, Z. Chen, C. Liu, M. Wu, C.-K. Kwoh, X. Li, and C. Guan, "An attention-based deep learning approach for sleep stage classification with single-channel eeg," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 29, pp. 809–818, 2021.
- [41] J. Pan, Y. Feng, P. Zhao, X. Zou, A. Hou, and X. Che, "Causalattennet: A fast and long-term-temporal network for automatic sleep staging with single-channel eeg," *IEEE Transactions on Instrumentation and Measurement*, vol. 73, pp. 1–13, 2024.