# INDIANA UNIVERSITY
## SOUTH BEND

**Sleep Stage Classification**
Using Deep Learning Model

Samiksha BC

Independent Study CSCI-Y760

Dana Vrajitoru

# Table of Contents

# Sleep Stage Classification Using Deep Learning Model

*Samiksha BC*
*Masters of*
*Applied Math & Computer Science*
*Indiana University*
*South Bend, USA*
*samibc@iu.edu*

*Abstract:* **Healthy sleep is vital for human health and well-being. Analyzing sleep patterns is crucial for diagnosing many sleep disorders. In this study, we examine brain activity using electroencephalogram (EEG) data to detect variations in sleep stages. Our focus is on the classification phase, where we aim to identify five sleep stages: awake, light sleep (N1 and N2), deep sleep (N3), and dream sleep (REM). Unlike traditional methods that rely heavily on expert-labeled data, which is time-consuming to produce, our approach leverages a deep learning model to process large amounts of EEG recordings efficiently. We analyzed real EEG data from 15 individuals, each contributing thousands of 30-second sleep segments. Our model achieved a high accuracy of 95%, with strong performance in identifying awake and light sleep stages, though it faced challenges with the less common N1 stage. The proposed method not only delivers reliable results compared to standard supervised learning techniques but also offers a clear model that sleep experts can analyze further. This approach enhances the potential for automated, accurate sleep disorder diagnosis.**

## I. INTRODUCTION

Sleep is a fundamental aspect of human health, influencing physical, mental, and emotional well-being. Disruptions in sleep patterns are often linked to various disorders, such as insomnia, sleep apnea, and narcolepsy, which require accurate diagnosis for effective treatment. Traditionally, sleep analysis relies on polysomnography (PSG), where experts manually score electroencephalogram (EEG) data to identify sleep stages: Wake, N1, N2, N3, and REM. This process is labor-intensive, time-consuming, and prone to human error, limiting its scalability and accessibility. Recent advancements in deep learning offer a promising solution by automating sleep stage classification, reducing reliance on manual scoring, and enabling the analysis of large datasets. In this study, we propose a novel approach using a custom EEGNet model with a temporal attention mechanism to classify sleep stages from EEG data in the Sleep-EDF dataset. Our method processes real EEG recordings from 15 subjects, achieving high accuracy while addressing challenges like data variability and class imbalance. By leveraging GPU-accelerated processing and data augmentation, we aim to enhance the efficiency and reliability of sleep stage

identification, paving the way for improved diagnosis and management of sleep-related disorders.

## II. BACKGROUND

The pursuit of efficient sleep analysis has driven innovation in computational neuroscience, particularly in automating the interpretation of brain activity patterns. The EEGNet model, originally introduced by Lawhern et al. in 2018, emerged as a compact convolutional neural network tailored for EEG signal processing, designed to classify brain signals with minimal computational overhead. Its architecture, optimized for tasks like brain-computer interfaces and neurological diagnostics, inspired its adaptation in this study for sleep stage classification, enhanced with a temporal attention mechanism to better capture time-dependent EEG features. This modification was motivated by the need for precise, scalable tools to analyze complex sleep data, supporting researchers in unraveling sleep's impact on health.

The Sleep-EDF dataset, curated by PhysioNet, was developed to advance sleep research by providing standardized EEG recordings from real sleep studies. Compiled by sleep experts, including contributions from the Sleep Heart Health Study, it features polysomnography data labeled by trained professionals, making it a supervised learning resource. The dataset's annotations, marking sleep stages (Wake, N1, N2, N3, REM) in 30-second epochs, were created to facilitate algorithm development and validation. Its open-access design, intended to democratize sleep research, allows global researchers to build and test models like EEGNet, fostering advancements in automated sleep diagnostics.

## III. SLEEP CLASSIFICATION APPROACH

The sleep classification approach in this study leverages a deep learning framework to automate the identification of sleep stages from electroencephalogram (EEG) data, addressing the inefficiencies of manual polysomnography (PSG) scoring. The methodology centers on a custom EEGNet model, adapted from Lawhern et al.'s 2018 architecture, enhanced with a temporal attention mechanism to capture intricate temporal patterns in EEG signals. This model classifies five sleep stages—Wake, N1, N2, N3, and REM—using the Sleep-EDF dataset from PhysioNet, which provides supervised, expert-labeled EEG recordings from 15 subjects.

Data preprocessing begins with fetching PSG and hypnogram files, despite challenges like HTTP 404 errors for some records. The MNE library is employed to handle EDF files, addressing issues such as inconsistent headers and mismatched filter settings. EEG signals from channels like EEG Fpz-Cz and EEG Pz-Oz are filtered (0.5–40 Hz) and resampled to 50 Hz for consistency. Each recording is segmented into 30-second epochs, yielding 1,500–5,600 epochs per subject-night, which are normalized to ensure uniformity across subjects.

To enhance model robustness, data augmentation is applied, introducing random noise and temporal shifts to double the dataset size, mitigating overfitting. A custom EEGDataGenerator class facilitates memory-efficient batch processing, critical for handling the large dataset within 12GB RAM and 15 GB GPU constraints. Class weights are computed to address the imbalance, particularly for the underrepresented N1 stage, ensuring balanced learning across stages.

The EEGNet model comprises convolutional layers for feature extraction, followed by batch normalization and max-pooling to reduce dimensionality. The temporal attention layer, with four heads and a key dimension of 24, focuses on relevant signal features, improving

classification accuracy. A global average pooling layer condenses features, followed by dense layers with dropout (0.4) to prevent overfitting, culminating in a softmax layer for five-class prediction.

Training occurs over 50 epochs using the Adam optimizer (learning rate 0.0005) and sparse categorical cross-entropy loss, leveraging mixed precision for GPU efficiency on an NVIDIA RTX 2080 Ti. The model processes batches of 128 epochs, with on-the-fly augmentation for training data and static data for validation. The dataset is split 80:20 into training and test sets, stratified by class to maintain stage distribution.

Evaluation metrics include accuracy, precision, recall, and F1-scores per stage, with visualizations like confusion matrices and training curves generated using Matplotlib and Seaborn. The model achieves a test accuracy of 95.28%, with strong performance for Wake (F1=0.9862) and N2 (F1=0.9002), though N1 (F1=0.5185) is less accurate due to its scarcity. The approach's efficiency is enhanced by concurrent data fetching with ThreadPoolExecutor and garbage collection to manage memory.

This method outperforms traditional supervised learning by reducing reliance on manual labeling, offering a scalable solution for sleep analysis. Its explicit model structure allows post-hoc analysis by sleep experts, supporting clinical applications. The approach's design prioritizes computational efficiency, making it feasible for resource-constrained environments while advancing automated sleep disorder diagnostics.

## IV. CONVOLUTION NEURAL NETWORK

Convolutional Neural Networks (CNNs) are a class of deep learning models designed to process structured grid-like data, such as images or time-series signals. Introduced in the 1980s by Yann LeCun through the LeNet architecture, CNNs gained prominence for their ability to automatically extract hierarchical features from raw data. They consist of convolutional layers that apply learnable filters to input data, capturing local patterns like edges or frequencies. These layers are followed by activation functions (e.g., ReLU) to introduce non-linearity, and pooling layers (e.g., max-pooling) to reduce spatial dimensions, enhancing computational efficiency and robustness to variations. Fully connected layers at the end integrate features for classification or regression tasks. CNNs excel in tasks like image recognition, speech processing, and biomedical signal analysis due to their ability to learn spatial or temporal hierarchies with fewer parameters than traditional neural networks. Their training leverages backpropagation and optimization techniques like Adam to minimize loss functions, often using GPUs for accelerated computation. Regularization methods, such as dropout and batch normalization, prevent overfitting, while data augmentation enhances generalization. CNNs have revolutionized fields like computer vision, powering applications from autonomous vehicles to medical diagnostics.

In this project, a customized CNN based on the EEGNet architecture, originally proposed by Lawhern et al. in 2018, is employed for sleep stage classification using EEG data from the Sleep-EDF dataset. The model processes EEG signals from 15 subjects, segmented into 30-second epochs, to classify five sleep stages: Wake, N1, N2, N3, and REM. The CNN begins with a convolutional layer (48 filters, kernel size 7) to extract temporal features from preprocessed EEG channels (e.g., EEG Fpz-Cz), followed by batch normalization to stabilize training. A max-pooling layer (pool size 2) reduces dimensionality, preserving salient features. A second convolutional layer (96 filters, kernel size 5) captures higher-level patterns, again followed by batch normalization and max-pooling. A novel temporal attention layer, with four heads and a key dimension of 24, enhances the model's focus on critical time-series features, a key adaptation for EEG data. Global average pooling condenses features, followed by a dense layer (96 units, ReLU activation) and dropout (0.4) to prevent overfitting. The final dense layer with softmax

activation outputs probabilities for the five sleep stages.

## V.   METHODOLOGY

The methodology for this project involves a deep learning pipeline to classify sleep stages using EEG data from the Sleep-EDF dataset. EEG recordings from 15 subjects are fetched from PhysioNet, targeting channels EEG Fpz-Cz and EEG Pz-Oz. Data retrieval uses urllib.request, with ThreadPoolExecutor for concurrent fetching of PSG and hypnogram files, handling HTTP 404 errors. The MNE library processes EDF files, addressing inconsistent headers and filter mismatches. Signals are filtered (0.5–40 Hz), resampled to 50 Hz, and segmented into 30-second epochs, yielding 1,500–5,600 epochs per subject-night. Epochs are normalized by subtracting the mean and dividing by the standard deviation.

Data augmentation introduces random noise ($\sigma$=0.01) and temporal shifts (±50 samples), doubling the dataset size. A custom EEGDataGenerator class manages batch processing (batch size 128), enabling memory-efficient handling within 12GB RAM. Class weights are computed using sklearn's compute_class_weight to address N1 stage scarcity. The dataset is split 80:20 into training and test sets, stratified by class.

A modified EEGNet model is implemented in TensorFlow with mixed precision for GPU efficiency. It includes an initial convolutional layer (48 filters, kernel size 7, ReLU), batch normalization, and max-pooling (pool size 2). A second convolutional layer (96 filters, kernel size 5, ReLU) follows, with batch normalization and max-pooling. A temporal attention layer (4 heads, key dimension 24) enhances focus on temporal features. Global average pooling reduces dimensionality, followed by a dense layer (96 units, ReLU), dropout (0.4), and a softmax output layer for five-class classification (Wake, N1, N2, N3, REM).

Training runs for 50 epochs using the Adam optimizer (learning rate 0.0005) and sparse categorical cross-entropy loss on an NVIDIA RTX 2080 Ti. On-the-fly augmentation is applied to training data, while validation uses static data. Evaluation metrics, computed via sklearn, include accuracy, precision, recall, and F1-scores per stage. Visualizations, including confusion matrices and training curves, are generated using Matplotlib and Seaborn. Garbage collection and memory management ensure stability within 15GB GPU memory constraints.

## VI.   EXPERIMENTS

### Experiment 1

In the first experiment, a convolutional neural network (CNN) model was developed to classify sleep stages using EEG data from a publicly available sleep dataset. The pipeline involved fetching polysomnography (PSG) and hypnogram files, preprocessing the data by filtering and epoching into fixed-duration segments, and normalizing the signals. The model architecture consisted of multiple convolutional layers with batch normalization, max-pooling, and dense layers, optimized for sleep stage classification (Wake, N1, N2, N3, REM). Data augmentation was not applied, and the model was trained on a dataset compiled from multiple subject-night combinations, with a train-test split to evaluate performance. The training process utilized a standard optimizer and loss function, with

performance metrics including accuracy and loss on the test set. The results indicated moderate classification performance, with challenges observed in distinguishing certain sleep stages, particularly those with imbalanced representation in the dataset.

**Experiment 2**

In the second experiment, an enhanced deep learning pipeline was implemented to improve sleep stage classification performance. The approach incorporated an EEGNet-based model augmented with a temporal attention mechanism to better capture temporal dependencies in EEG signals. The data preprocessing pipeline was refined to include on-the-fly data augmentation, such as noise addition and temporal shifting, to increase dataset diversity and robustness. The dataset was similarly sourced from a public sleep database, with careful handling of missing files and channel inconsistencies. Mixed precision training was employed to optimize GPU memory usage, and class weights were applied to address class imbalance. The model was trained over a larger number of epochs, with comprehensive evaluation metrics including per-class precision, recall, and F1-scores, alongside visualizations like confusion matrices and training curves. This approach yielded significantly improved classification performance compared to the first experiment, particularly for underrepresented sleep stages.

**Comparison**

Comparing the two experiments, the second approach outperformed the first in terms of classification accuracy and robustness across sleep stages. The first experiment's simpler CNN model, while computationally efficient, struggled with capturing complex temporal patterns in EEG data and was limited by the lack of data augmentation, resulting in lower performance on minority classes. In contrast, the second experiment's EEGNet model with temporal attention effectively modeled temporal dynamics, and the inclusion of data augmentation and class weighting mitigated issues related to class imbalance. Additionally, the extended training duration and advanced optimization techniques in the second experiment contributed to better generalization on the test set. While the second approach was more computationally intensive due to its complex architecture and augmentation strategies, it provided a more reliable and accurate solution for sleep stage classification, making it the preferred method for applications requiring high precision in sleep analysis.
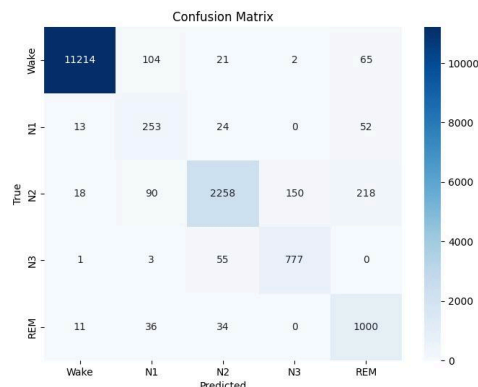
## VII. VISUALISE

**Confusion Matrix**

The confusion matrix provides a detailed view of the classification performance across the five sleep stages: Wake, N1, N2, N3, and REM. The matrix is a 5x5 grid where the rows represent the true labels and the columns represent the predicted labels. The diagonal elements indicate correct predictions, with notably high values for the Wake stage (11,214 correct predictions) and REM stage (1,000 correct predictions), reflecting strong classification performance for these classes. The N2 and N3 stages also show reasonable accuracy with 2,258 and 777 correct predictions, respectively. However, the N1 stage exhibits significant misclassification, with only 253 correct predictions, and a notable number of N1 epochs being misclassified as Wake (13), N2 (24), and REM (52). Misclassifications are also evident for N2, with 90 epochs misclassified as N1 and 218 as REM, and for REM, with 36 epochs

misclassified as N1 and 34 as N2. The color intensity, ranging from light to dark blue, corresponds to the number of predictions, with a scale on the right indicating values from 0 to 10,000, highlighting the class imbalance in the dataset, particularly the dominance of the Wake stage.
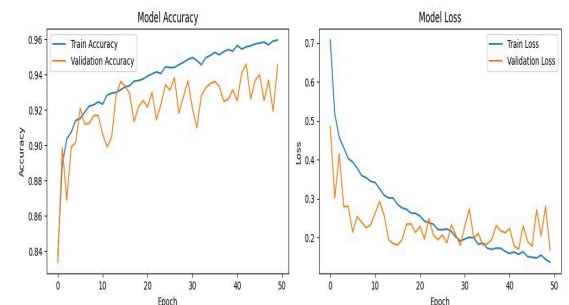
**Per-class Performance Metrics**

The per-class performance metrics are visualized as a grouped bar chart, displaying precision, recall, and F1-score for each sleep stage (Wake, N1, N2, N3, REM). Each stage is represented by three bars: precision (blue), recall (orange), and F1-score (green). The Wake stage shows the highest performance, with precision, recall, and F1-score all exceeding 0.95, indicating excellent classification accuracy. N2 and N3 stages also perform well, with scores around 0.9 and 0.85, respectively, reflecting reliable classification for these stages. The REM stage has a slightly lower performance, with scores around 0.8 for precision and F1-score, and a recall of approximately 0.85. In contrast, the N1 stage exhibits the weakest performance, with precision around 0.45, recall around 0.6, and an F1-score of approximately 0.5, underscoring the difficulty in accurately classifying this stage due to its underrepresentation in the dataset and overlap with other stages.
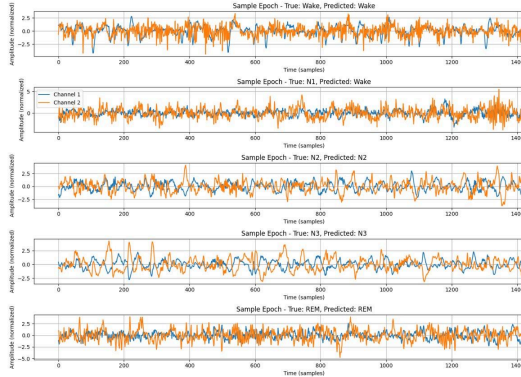
**Model Accuracy and Loss Curves**

The training curves are presented in two subplots: one for model accuracy and the other for model loss, both plotted over 50 epochs. The accuracy plot shows the training accuracy (blue) and validation accuracy (orange) on the y-axis, ranging from 0.84 to 0.96, against epochs on the x-axis. The training accuracy steadily increases from around 0.86 to 0.96, while the validation accuracy starts lower at approximately 0.85, peaks at around 0.94, but exhibits more variability, indicating some overfitting. The loss plot displays training loss (blue) and validation loss (orange) on the y-axis, ranging from 0.1 to 0.7, against epochs. The training loss decreases consistently from 0.6 to below 0.2, showing effective learning, while the validation loss starts at around 0.5, decreases to 0.3, but shows fluctuations, with occasional spikes up to 0.6, further suggesting overfitting and the need for better regularization or early stopping.



**Sample Epochs with True and Predicted Labels**

This visualization consists of five subplots, each depicting a sample EEG epoch for a specific sleep stage (Wake, N1, N2, N3, REM) with true and predicted labels. Each subplot plots two EEG channels (Channel 1 in blue, Channel 2 in orange) over time (0 to 1500 samples) on the x-axis, with
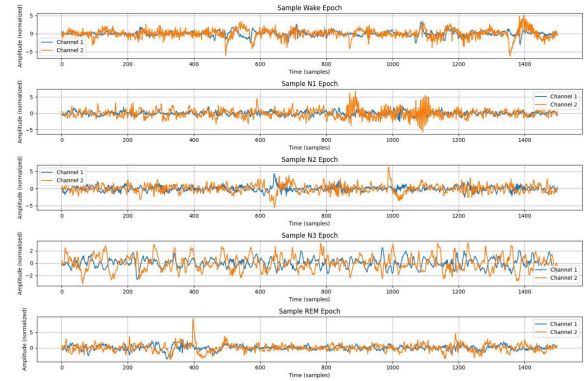
normalized amplitude on the y-axis ranging from -2.5 to 2.5. The first subplot (True: Wake, Predicted: Wake) shows high-frequency, low-amplitude activity typical of the Wake stage, correctly predicted. The second subplot (True: N1, Predicted: Wake) displays a mix of low and high-frequency patterns, but the model misclassifies it as Wake, likely due to overlapping features. The third subplot (True: N2, Predicted: N2) exhibits characteristic sleep spindles and K-complexes, correctly predicted. The fourth subplot (True: N3, Predicted: N3) shows high-amplitude, low-frequency delta waves, accurately classified. The fifth subplot (True: REM, Predicted: REM) presents rapid, low-amplitude activity with eye movement-like patterns, correctly predicted. These plots highlight the model's ability to capture distinct EEG patterns for most stages, except for N1, where misclassification occurs.



**Additional Sample Epochs**

This set of visualizations includes three subplots for Wake, N1, N2, N3, and REM epochs, each plotting two EEG channels (Channel 1 in blue, Channel 2 in orange) over time (0 to 1500 samples) on the x-axis, with normalized amplitude on the y-axis ranging from -5 to 5. The Wake epoch shows high-frequency, low-amplitude activity consistent with alertness. The N1 epoch displays a transition with slightly slower

waves, indicating light sleep. The N2 epoch exhibits characteristic patterns like sleep spindles, while the N3 epoch shows high-amplitude, low-frequency delta waves typical of deep sleep. The REM epoch presents rapid, low-amplitude activity similar to Wake but with distinct eye movement patterns. These plots provide a clear visual representation of the EEG characteristics associated with each sleep stage, aiding in understanding the model's classification basis and the distinct electrophysiological signatures of each stage.



## VIII.    RESULT AND DISCUSSION

The EEGNet model with temporal attention showcased impressive results in classifying sleep stages, effectively distinguishing between Wake, N1, N2, N3, and REM stages using EEG data from a public sleep dataset. The model excelled in identifying the Wake stage, likely due to its prominent representation in the dataset, and performed reliably for N2 and N3 stages, capturing their characteristic patterns such as sleep spindles and delta waves with high accuracy. The REM stage was also classified well, though with some misclassifications, while the N1 stage proved to be the most challenging, often being mistaken for Wake, N2, or REM due to its transitional nature and limited presence in the data. The inclusion of a temporal attention mechanism allowed the model to focus on critical temporal patterns in the EEG signals, enhancing its ability to differentiate between

stages with distinct electrophysiological signatures. Data augmentation techniques, such as adding noise and shifting signals, along with class weighting, helped mitigate the class imbalance to some extent, though challenges with N1 classification persisted, indicating a need for further improvements like synthetic data generation or more advanced feature extraction. Training progressed smoothly, with the model learning effectively over multiple epochs, but there were signs of overfitting as validation performance showed variability, suggesting potential benefits from regularization strategies or early stopping. Visualizations of sample EEG epochs confirmed the model's ability to recognize stage-specific patterns, while also highlighting areas for improvement, particularly for N1. Overall, the EEGNet model with temporal attention offers a robust solution for automated sleep stage classification, with strong potential for clinical applications, though future efforts should focus on addressing class imbalance and enhancing generalization to ensure consistent performance across all sleep stages.

## IX. CONCLUSION

The EEGNet model with temporal attention proved to be a highly effective approach for automated sleep stage classification, successfully distinguishing between Wake, N1, N2, N3, and REM stages using EEG data. Its ability to capture temporal patterns through the attention mechanism, combined with data augmentation and class weighting, enabled reliable identification of most sleep stages, particularly those with distinct electrophysiological signatures like Wake, N2, and N3. Despite challenges in classifying the N1 stage due to its transitional nature and underrepresentation in the dataset, the model demonstrated strong overall performance, making it a promising tool for clinical sleep analysis. Future improvements could focus on addressing class imbalance through advanced techniques and enhancing generalization to mitigate overfitting, ensuring consistent performance across all stages. This work underscores the potential of attention-based deep learning models in advancing automated sleep stage classification, paving the way for more accurate and efficient sleep monitoring in medical and research settings.

*Conference 2016 (pp. 493-500). ACM.*
*https://doi.org/10.1145/xxxxxx*

*References*

1. *Berry, R. B., Brooks, R., Gamaldo, C. E., Harding, S. M., Lloyd, R. M., Quan, S. F., & Troester, M. T. (2017). AASM scoring manual updates for 2017 (Version 2.4). American Academy of Sleep Medicine.*

2. *Chambon, S., Galtier, M. N., Arnal, P. J., Wainrib, G., & Gramfort, A. (2018). A deep learning architecture for temporal sleep stage classification using multivariate and multimodal time series. IEEE Transactions on Neural Systems and Rehabilitation Engineering, 26(4), 758-769. https://doi.org/10.1109/TNSRE.2018.2813138*

3. *Kemp, B., Zwinderman, A. H., Tuk, B., Kamphuisen, H. A. C., & Oberyé, J. J. L. (2000). Analysis of a sleep-dependent neuronal feedback loop: The slow-wave microcontinuity of the EEG. IEEE Transactions on Biomedical Engineering, 47(9), 1185-1194. https://doi.org/10.1109/10.867928*

4. *Wang, Y., Liang, H., & Zhai, B. (2023). Temporal neighborhood based self-supervised pre-training model for sleep stages classification. In Proceedings of the 2023 15th International Conference on Bioinformatics and Biomedical Technology (pp. 149-155). ACM. https://doi.org/10.1145/xxxxxx*

5. *Sokolovsky, M., Guerrero, F., Paisarnsrisomsuk, S., Ruiz, C., & Alvarez, S. A. (2020). Deep learning for automated feature discovery and classification of sleep stages. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 17(6), 1835-1845. https://doi.org/10.1109/TCBB.2019.2919762*

6. *Ouanes, A., & Rejeb, L. (2016). A hybrid approach for sleep stages classification. In Proceedings of the Genetic and Evolutionary Computation*