

Explainable Sleep Stage Classification with Multimodal Electrophysiology Time-series*

Charles A. Ellis, Rongen Zhang, Darwin A. Carbajal, Robyn L. Miller, Vince D. Calhoun, May D. Wang

Abstract— Many automated sleep staging studies have used deep learning approaches, and a growing number of them have used multimodal data to improve their classification performance. However, few studies using multimodal data have provided model explainability. Some have used traditional ablation approaches that “zero out” a modality. However, the samples that result from this ablation are unlikely to be found in real electroencephalography (EEG) data, which could adversely affect the importance estimates that result. Here, we train a convolutional neural network for sleep stage classification with EEG, electrooculograms (EOG), and electromyograms (EMG) and propose an ablation approach that replaces each modality with values that approximate the line-related noise commonly found in electrophysiology data. The relative importance that we identify for each modality is consistent with sleep staging guidelines, with EEG being important for most sleep stages and EOG being important for Rapid Eye Movement (REM) and non-REM stages. EMG showed low relative importance across classes. A comparison of our approach with a “zero out” ablation approach indicates that while the importance results are consistent for the most part, our method accentuates the importance of modalities to the model for the classification of some stages like REM ($p < 0.05$). These results suggest that a careful, domain-specific selection of an ablation approach may provide a clearer indicator of modality importance. Further, this study provides guidance for future research on using explainability methods with multimodal electrophysiology data.

Clinical Relevance— While explainability is helpful for clinical machine learning classifiers, it is important to consider how explainability methods interact with clinical data, a domain for which they were not originally designed.

I. INTRODUCTION

Many methods have been developed for automated sleep staging in recent years. Most use electroencephalograms (EEG) [1] or electrooculograms (EOG) [2], and only a few utilize multimodal data [3], [4]. Clinicians typically use multimodal data when scoring sleep stages. As such, the use

of multimodal data, like EEG, EOG, and electromyograms (EMG) could provide insights that could not be gained from unimodal data alone. Multimodal data facilitates more intricate recognition of human activity [4]. Moreover, while many studies use deep learning (DL) for automated sleep staging, most do not give insight into the inner mechanisms of their classifiers [5]. The black-box nature of DL models is problematic for clinical implementation because they are difficult for clinicians to interpret. Of the papers that offer explainability, most involve EEG [6]–[8]. Studies on multimodal data, with a few exceptions [5], [6], do not use explainability methods or provided insight into the importance of the modalities that they analyze. In this study, we present a novel explainability approach that enables us to identify the importance of different modalities to a classifier and that is better suited to explaining multimodal electrophysiology (EP) classifiers than preexisting approaches.

Several studies use explainable artificial intelligence methods to examine modality importance [9], [10]. These studies use methods like layer-wise relevance propagation (LRP) and ablation. Ablation involves the removal of each modality and the calculation of the effect that its removal has upon the classifier performance. Like similar methods that perturb data, ablation methods can create samples that are outside the data distribution upon which the classifier is trained and potentially give an inadequate explanation [11]. While ablation approaches are simple and intuitive, it is important to consider the potential problems that can arise when they are applied in domains, like electrophysiology classification, for which they were not originally designed. Existing multimodality studies that use ablation replace the values of each modality with zeroes [9], [10]. This “zeroing out” of a modality may not force a sample outside the data distribution if an appropriate data normalization method (i.e., z-scoring) is used. However, the resulting samples are, nevertheless, unrealistic and do not align with how real-life samples appear.

In this study, we train a 1-dimensional (1D) convolutional neural network (CNN) for automated sleep staging with

*Funding for this work is from NIH grant R01EB006841.

C. A. Ellis is with the Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology and Emory University, Atlanta, GA 30332 USA. (corresponding author, e-mail: cae67@gatech.edu).

R. Zhang is with the Department of Computer Information Systems, Georgia State University, Atlanta, GA 30303 USA. (e-mail: rzhang6@gsu.edu)

D. A. Carbajal is with the Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA 30332 USA. (email: darwin.carbajal@gatech.edu)

R. L. Miller is with the Tri-institutional Center for Translational Research in Neuroimaging and Data Science: Georgia State University, Georgia Institute of Technology, Emory University, Atlanta, GA 30303 USA. (e-mail: robyn.l.miller@gmail.com)

V. D. Calhoun is with the Tri-institutional Center for Translational Research in Neuroimaging and Data Science: Georgia State University, Georgia Institute of Technology, Emory University, Atlanta, GA 30303 USA. (e-mail: vcalhoun@gsu.edu)

M. D. Wang is with the Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology and Emory University, Atlanta, GA 30332 USA. (email: maywang@gatech.edu)

multimodal electrophysiology data. Then to address the previously described problem of explainability, we present a novel ablation-based explainability method that finds the relative importance of each modality to the classification of each sleep stage. Our ablation approach seeks to ablate each modality in a way that creates new samples that (1) a classifier might encounter during a clinical implementation and (2) that would not be as distinct from the original data as “zeroed out” samples. We do this by replacing each modality with a sinusoid and Gaussian noise that mimics the line-related noise that is commonly found in electrophysiology recordings. To examine how our method relates to existing approaches, we compare the results of our explainability method with that of the traditional “zeroing out” ablation approach.

II. METHODS

Here we describe our study approach. We train a 1D CNN for sleep stage classification with multimodal data and output explanations for insight into the modalities critical to classifying each stage using our novel ablation method. We then implement an existing ablation method and use statistical tests to compare the two methods.

A. Description of Data

We use sleep telemetry data from the PhysioNet [12] Sleep-EDF Database [13]. The data was not collected specifically for our study, and no one on our team had access to subject identifiers. As such, our study was not considered human subjects research and did not require Institutional Review Board approval. The dataset consists of 44 full night (approximately 9 hour) recordings collected from 22 subjects with primary onset insomnia. Each subject has two recordings: one following placebo administration and one following temazepam administration. The data has three modalities: EEG, EOG, and EMG. All modalities have a sampling rate of 100 Hertz (Hz). For EEG, we use the FPz-Cz electrode. A marker recorded at 1 Hz indicates whether the sleep telemetry system operated correctly, and a polysomnogram consisting of Awake, Movement, rapid eye movement (REM), Non-REM 1 (NREM1), NREM2, NREM3, and NREM4 is also included.

B. Description of Preprocessing

We divide the data into non-overlapping, 30-second segments and extract labels from the polysomnograms. NREM3 and NREM4 stages are combined into a single NREM3 stage [14] and Movement samples are discarded. We discard all samples that coincide with a recording error. We apply z-score normalization to each modality within each

recording. After segmentation, the dataset has 42,218 samples. Approximately, 9.97% (4,213 samples), 8.53% (3,603 samples), 46.8% (19,755 samples), 14.92% (6,298 samples), and 19.78% (8,349 samples) belong to the Awake, NREM1, NREM2, NREM3, and REM classes, respectively.

C. Convolutional Neural Network

We adapt a CNN architecture that was first developed for EEG sleep stage classification [15]. Similar architectures have been used in previous studies [8]. The architecture is shown in Figure 1, and we implement it in Tensorflow and Keras. We use a 1D-CNN because CNNs extract relevant features from time-series, and sleep EP data has many relevant features (e.g., EEG frequency bands, EMG spikes). When training the model, we use 10-fold cross-validation in which 17, 2, and 3 subjects in each fold are randomly assigned to training, validation, and test groups, respectively. While training the classifier, we use categorical cross entropy loss and weight the loss for each class to account for class imbalances. We use the Adam optimizer [16] with an adaptive learning rate that decreases after every 5 epochs with no increase in validation accuracy. The optimizer has an initial learning rate of 0.001. During testing of each fold, we use the model weights from the epoch with the best validation accuracy. To account for class imbalances when assessing test performance, we calculate the precision, recall, and F1 score for each class in each fold. We then calculate their mean and standard deviation (SD) across all folds.

D. Ablation-based Global Explainability

We apply an ablation approach for insight into the importance of each modality to the identification of each sleep stage. Instead of zeroing out each modality, we replace them with values that might be expected in EP data. Line-related noise often appears in EP data at around 50 Hz or 60 Hz as a result of the presence of power lines, lights, and other electronics near recording devices, and when an electrode is not working properly, it is common to find only line-related noise in that particular channel. For a sampling rate of 100 Hz, aliased 60 Hz noise should appear at around 40 Hz. As such, for our study, we replace each modality with a combination of a 40 Hz sinusoid with an amplitude of 0.1 and Gaussian noise with a mean of 0 and SD of 0.1. Before ablation, we measure the weighted F1 score across all classes and the F1 scores for each individual class. We then ablate each modality and calculate the resulting change in performance (Original F1 – New F1). We perform the

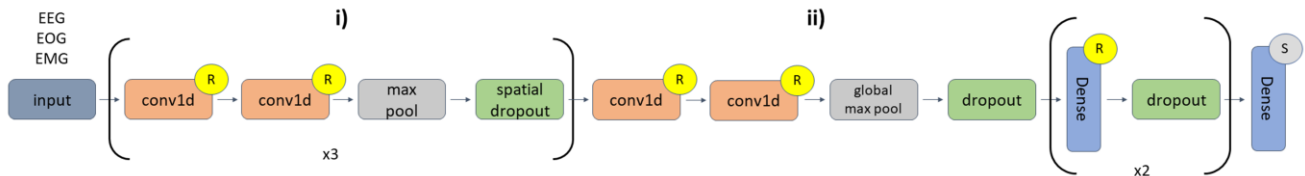


Figure 1. CNN Architecture. In the repeated layer i) of the diagram there are 6 1D convolutional (conv1d) layers. The first two conv1d layers have 16 filters with a kernel size of 5 followed by a max pooling layer with a pool size of 2 and a spatial dropout layer with a rate of 0.01. The second two conv1d layers (in i) had 32 filters with a kernel size of 3, followed by a max pooling layer with a pool size of 2 and a spatial dropout layer with a rate of 0.01. The third pair of conv1d layers (in i) have 32 filters with a kernel size of 3 followed by max pooling with a pool size of 2 and spatial dropout with a rate of 0.01. The last two conv1d layers (in ii) have 256 filters with a filter size of 3 followed by global max pooling and dropout with a rate of 0.01. The first dense layer has 64 nodes with dropout layer with a rate of 0.1. The second dense layer has 64 nodes with a dropout layer with a rate of 0.05. The last dense layer has 5 nodes. Layers with an “R” or an “S” indicate that they are followed by a ReLU or Softmax activation function, respectively.

TABLE I. CLASSIFICATION PERFORMANCE RESULTS

	Awake	NREM1	NREM2	NREM3	REM
Precision	72.25±07.12	36.20±03.98	79.35±03.92	56.78±18.35	69.04±07.14
Recall	70.90±07.02	46.28±13.52	68.71±08.51	78.22±10.24	63.26±06.69
F1	71.25±05.15	39.86±07.19	73.28±04.76	64.15±15.25	65.92±06.28

ablation for each fold individually. For a comparison, we use a “zeroing-out” ablation approach [9] and perform two-tailed t-tests to compare the change in F1 associated with each class and modality for the two methods.

III. RESULTS AND DISCUSSION

In this section, we discuss the model performance and the insights gained by comparing the ablation methods.

A. Model Performance Results

Table 1 shows the test performance of the model across all folds. The classification of NREM1 samples obtains the lowest level of performance across all metrics, which is unsurprising given that the NREM1 stage is the smallest class. Although the Awake and NREM1 classes have a comparable number of samples, the model obtains much higher performance for Awake. Additionally, the classifier obtains higher precision for the Awake class than all other classes except for NREM2, higher recall than all other classes except for NREM3, and a higher F1 score than all other classes except for NREM2. This makes sense given that EEG, EOG, and EMG Awake activity is very different from NREM activity and that EMG Awake activity is very different from EMG REM activity [14]. Additionally, the model obtains highest precision and F1 scores for the NREM2 class, which is reasonable given that nearly half of the dataset is NREM2.

B. Ablation-based Explainability Results

Figure 2 shows the explainability results for each method, along with the results of the statistical analysis. Panel A shows

the change in F1 across all classes, and Panels B through F show the change in F1 for each individual class. For the line noise-related ablation analysis with the F1 score calculated for all classes, EEG is the most important modality by far, with a median reduction in F1 of nearly 60% following the ablation of EEG. This change in F1 for all classes is skewed by the NREM2 and NREM3 classes. For the NREM2, NREM 3, and Awake classes, EEG is by far the most important modality while EOG and EMG have relatively little effect. EOG and EMG have larger effects upon the F1 score for the NREM1 and REM classes. For the NREM1 class, EEG and EOG have comparable effects upon the F1 score, though EEG has a slightly higher effect. Interestingly, the ablation of EMG for the NREM1 class seems to have a beneficial effect, increasing the F1 score by as much as 6-7%. For the REM class, both EEG and EOG have a significant effect upon the F1 score, though EEG has a markedly larger effect. Additionally, for the REM class, EMG has an effect upon the F1 score as high as 15-16%. Given that our architecture was originally designed for EEG classification, it is possible that it did not effectively extract EMG features. This could explain the relatively low importance of EMG to the model.

The results of zeroing out each modality are highly similar to the results of applying line-related noise ablation. Most pairs of F1 ablation scores do not have a significant difference between them. However, when including all classes, EMG importance is significantly different between ablation methods ($p < 0.01$). For individual classes, the EEG of the Awake class ($p < 0.05$) and the EMG of the NREM2 class ($p < 0.05$) have

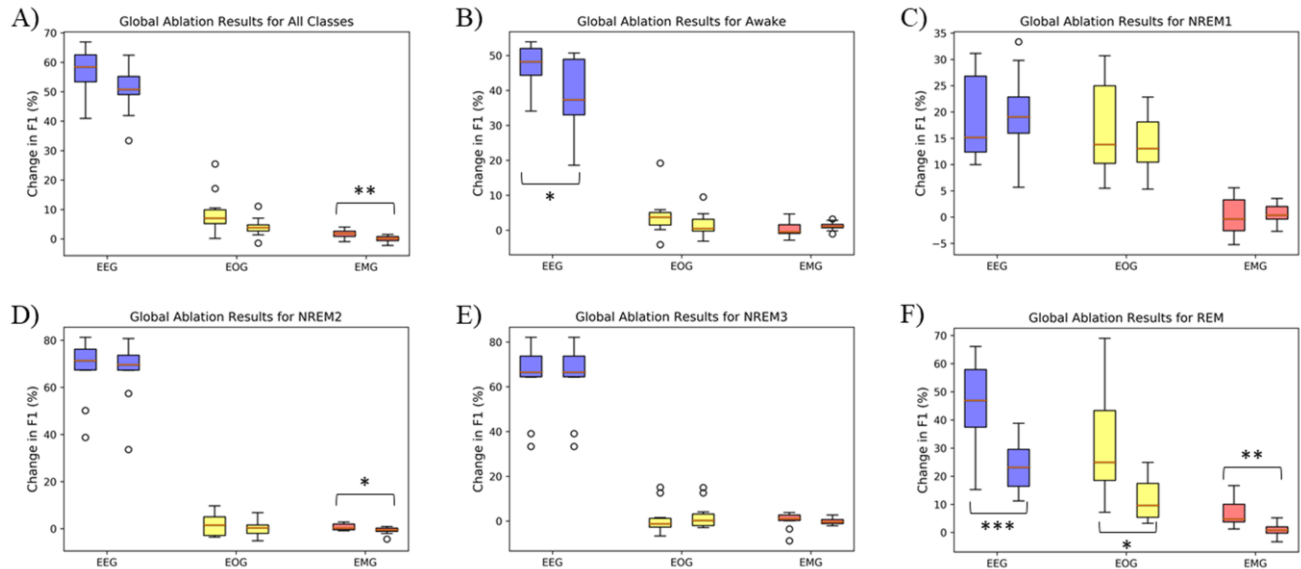


Figure 2. Results for Explainability Methods and Statistical Analysis. Panel A shows the explainability results across all classes, and Panels B through F show the results for individual classes. Blue, yellow, and red bars show the results for EEG, EOG, and EMG ablation, respectively. The leftmost of each pair of boxes shows the results for our ablation approach, and the rightmost shows the results for zeroing out each modality. The y-axes show the percent change in F1 following ablation, where positive and negative changes indicate decreases and increases in the F1 score, respectively. Some pairs of boxes are accompanied by *, **, or ***, which correspond to a two-tailed t-test p-value less than 0.05, 0.01, and 0.001, respectively.

significant results. Interestingly, all modalities have significant differences for the REM category ($p < 0.05$). In general, when a significant difference exists between importance assigned to a modality by the two explainability methods, our ablation method seems to give greater importance to the modality. This may indicate that our approach accentuates the importance of some modalities, which could be attributed to its use of ablation values that are more similar to real data.

The results of both methods generally fit with sleep scoring guidelines [14]. It is expected that EEG would be most important given that EEG varies significantly between NREM, Awake, and REM stages [14]. While relying upon EEG would enable the classifier to obtain good differentiation between Awake/REM and NREM, relying upon EEG may not help with classification between Awake and REM to the same degree. It would be logical for EMG to be important for classifying Awake and REM samples, as more movement might be expected in Awake than in NREM and as REM EMG activity would be much less than Awake and NREM activity [14]. It is interesting that EEG and EOG have comparable importance to NREM1. Given that the classifier obtains lowest performance for the NREM1 class, it may inappropriately rely upon EOG for identifying NREM1 samples.

D. Future Work

Examining model architectures that might better extract EMG features could be beneficial. While our explainability results fit with sleep scoring guidelines, suggesting the broader generalizability of the classifier, our explainability could potentially be improved. We seek to provide an alternative to traditional ablation approaches that zero out a feature by instead replacing the modalities with values similar to artifacts that naturally occur in EP recordings. When we approximate line-related noise, we have three parameters: the amplitude of the sinusoid and the mean and SD of the Gaussian noise. It is possible that using other parameter values might improve explanations. Also, exploring explainability methods that do not require modifying samples could improve explanations.

A. Conclusion

In this study, we train a classifier for multimodal sleep stage classification. We further propose a novel ablation-based approach that provides more realistic ablated samples than methods that simply zero out a particular modality. A comparison of our method with the traditional ablation approach indicates that the importance values are comparable with the exception of a few instances in which our method seems to accentuate the importance of some modalities. More broadly, our work has implications for ablation-based explainability in other data types. Specifically, it suggests that more careful consideration of how features are ablated in light of their domain may provide increases in importance metrics and clarify the relative importance of features.

ACKNOWLEDGMENT

We thank Felipe Giuste and Wenqi Shi for their advice and assistance with computational resources. We thank Mohammad Sendi for helping edit the paper.

REFERENCES

- [1] A. Sors, S. Bonnet, S. Mirek, L. Vercueil, and J. F. Payen, "A convolutional neural network for sleep stage scoring from raw single-channel EEG," *Biomed. Signal Process. Control*, vol. 42, pp. 107–114, 2018, doi: 10.1016/j.bspc.2017.12.001.
- [2] M. M. Rahman, M. I. H. Bhuiyan, and A. R. Hassan, "Sleep stage classification using single-channel EOG," *Comput. Biol. Med.*, vol. 102, no. June, pp. 211–220, 2018, doi: 10.1016/j.compbimed.2018.08.022.
- [3] S. Chambon, M. N. Galtier, P. J. Arnal, G. Wainrib, and A. Gramfort, "A deep learning architecture for temporal sleep stage classification using multivariate and multimodal time series," *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 26, no. 4, pp. 758–769, 2018.
- [4] B. Zhai, I. Perez-Pozuelo, E. A. D. Clifton, J. Palotti, and Y. Guan, "Making Sense of Sleep: Multimodal Sleep Stage Classification in a Large, Diverse Population Using Movement and Cardiac Sensing," *Proc. ACM Interactive, Mobile, Wearable Ubiquitous Technol.*, vol. 4, no. 2, 2020, doi: 10.1145/3397325.
- [5] O. Tsinalis, P. M. Matthews, and Y. Guo, "Automatic Sleep Stage Scoring Using Time-Frequency Analysis and Stacked Sparse Autoencoders," *Ann. Biomed. Eng.*, vol. 44, no. 5, pp. 1587–1597, 2016, doi: 10.1007/s10439-015-1444-y.
- [6] S. Mousavi, F. Afghah, and U. Rajendra Acharya, "SleepEEGNet: Automated Sleep Stage Scoring with Sequence to Sequence Deep Learning Approach," *arXiv*, pp. 1–15, 2019, doi: 10.13026/C2C30J.
- [7] A. Vilamala, K. H. Madsen, and L. K. Hansen, "Deep convolutional neural networks for interpretable analysis of EEG sleep stage scoring," *IEEE Int. Work. Mach. Learn. Signal Process. MLSP*, vol. 2017-Sept, no. 659860, pp. 1–6, 2017, doi: 10.1109/MLSP.2017.8168133.
- [8] O. Tsinalis, P. M. Matthews, Y. Guo, and S. Zafeiriou, "Automatic Sleep Stage Scoring with Single-Channel EEG Using Convolutional Neural Networks," *arXiv*, 2016, [Online]. Available: <http://arxiv.org/abs/1610.01683>.
- [9] S. Pathak, C. Lu, S. B. Nagaraj, M. van Putten, and C. Seifert, "STQS: Interpretable multi-modal Spatial-Temporal-sequential model for automatic Sleep scoring," *Artif. Intell. Med.*, vol. 114, no. January, p. 102038, 2021, doi: 10.1016/j.artmed.2021.102038.
- [10] J. Lin, S. Pan, C. S. Lee, and S. Oviatt, "An Explainable Deep Fusion Network for Affect Recognition Using Physiological Signals," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 2069–2072, doi: <https://doi.org/10.1145/3357384.3358160>.
- [11] C. Molnar, *Interpretable Machine Learning A Guide for Making Black Box Models Explainable*, 2018th-08–14th ed. Lean Pub, 2018.
- [12] G. AL *et al.*, "PhysioBank, PhysioToolkit, and PhysioNet: Components of a New Research Resource for Complex Physiologic Signals," *Circulation*, vol. 101, no. 23, pp. e215–e220, 2000, [Online]. Available: <http://circ.ahajournals.org/content/101/23/e215.full>.
- [13] B. Kemp, A. H. Zwiderman, B. Tuk, H. A. C. Kamphuisen, and J. J. L. Obery, "Analysis of a sleep-dependent neuronal feedback loop: the slow-wave microcontinuity of the EEG," *IEEE Trans. Biomed. Eng.*, vol. 47, no. 9, pp. 1185–1194, 2000, doi: 10.1109/10.867928.
- [14] C. Iber, S. Ancoli-Israel, A. L. Chesson, and S. F. Quan, "The AASM Manual for Scoring of Sleep and Associated Events: Rules, Terminology, and Technical Specifications." 2007.
- [15] M. Younes, "CVxTz/EEG_classification: v1.0," 2020. https://github.com/CVxTz/EEG_classification (accessed Jan. 05, 2021).
- [16] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv Prepr. arXiv1412.6980*, 2014.