

# Multi-BK-Net: Multi-Branch Multi-Kernel Convolutional Neural Networks for Clinical EEG Analysis

Anonymous authors

Paper under double-blind review

## Abstract

Classifying an electroencephalography (EEG) recording as pathological or non-pathological is an important first step in diagnosing and managing neurological diseases and disorders. As manual EEG classification is costly, time-consuming and requires highly trained experts, deep learning methods for automated classification of general EEG pathology offer a promising option to assist clinicians in screening EEGs. Convolutional neural networks (CNNs) are well-suited for classifying pathological EEG signals due to their ability to perform end-to-end learning. In practice, however, current CNN solutions suffer from limited generalisation due to I) a single-scale network design that cannot fully capture the high intra- and inter-subject variability of the EEG signal, the diversity of the data, and the heterogeneity of pathological EEG patterns; and II) the small size and limited diversity of the dataset commonly used to train and evaluate the networks. These challenges result in a low sensitivity score and a performance drop on other datasets, further hindering their reliability for real-world applications. Here, we propose a novel multi-branch, multi-scale CNN called Multi-BK-Net (Multi-Branch Multi-Kernel Network), comprising five parallel branches that incorporate temporal convolution, spatial convolution, and pooling layers, with temporal kernel sizes defined based on five clinically relevant frequency bands within its first block. Evaluation is based on two public datasets with predefined test sets: the Temple University Hospital (TUH) Abnormal EEG Corpus and the TUH Abnormal Expansion Balanced EEG Corpus. Our Multi-BK-Net outperforms five baseline architectures and state-of-the-art end-to-end approaches in terms of accuracy and sensitivity on these datasets, setting a new benchmark. Furthermore, ablation experiments highlight the effectiveness of the multi-branch, multi-scale input block of the Multi-BK-Net. Overall, our approach demonstrates the effectiveness of multi-branch, multi-scale CNNs in accurately and reliably classifying general EEG pathology, while being more effective at handling data heterogeneity, and constitutes a next step towards deep end-to-end classification of general EEG pathology.

## 1 Introduction

Electroencephalography (EEG) is a non-invasive method of measuring and recording electrical activity in the brain. EEG has high temporal resolution, low equipment cost and is highly sensitive to dynamic changes in neural signals. Due to its high efficiency and usability, EEG is most commonly used in clinical practice for the diagnosis and management of various neurological conditions (Lopez et al., 2015; Zhang et al., 2023b). For this purpose, a preliminary important step in clinical practice is to classify an EEG recording as non-pathological or pathological (Brogger et al., 2018; Lopez et al., 2015). To this end, human EEG experts visually analyse recordings from long-term monitoring or multiple short sessions, which is a tedious and time-consuming process (Brogger et al., 2018). They also consider various additional patient-related factors, such as medical history, age, or medication (Beuchat et al., 2021; Nayak & Anilkumar, 2020). Furthermore, it requires years of training to achieve a thorough understanding of pathological EEG patterns and to distinguish them from normal EEG, normal benign variants, and artefacts (Amin et al., 2023; Emmady & Anilkumar, 2023; Hoppe, 2018). Thus, these challenges result in inter-rater variability and diagnostic errors. For the task of classifying EEG recordings as pathological or non-pathological, previous research has reported

an inter-rater agreement of 86–88% between two neurologists (Houfek & Ellingson, 1959; Rose et al., 1973), while Beuchat et al. (2021) has found even lower inter-rater agreements of 82-86% between multiple EEG technologists and neurologists.

In this sense, the introduction of automated EEG classification, or at least some level of automation such as clinical decision support systems, has the potential to improve or at least accelerate the EEG classification process and thereby enhance the quality of patient care. To this end, the use of deep learning approaches for EEG classification has received increasing attention in recent years (for a detailed review, see Amrani et al., 2021; Craik et al., 2019; Faust et al., 2018; Rahman et al., 2024; Roy et al., 2019b; Praveena et al., 2022) and has also spurred their application to the task of classifying general EEG pathology (Schirrmeyer et al., 2017b). In this regard, research has shown that convolutional neural networks (CNNs) have been quite effective for general EEG pathology classification due to their ability to extract relevant feature representations directly from raw or minimally preprocessed EEG data (Darvishi-Bayazi et al., 2023; Gemein et al., 2020; Khan et al., 2022; Kiessner et al., 2023; 2024; Van Leeuwen et al., 2019; Western et al., 2021, Appendix A.1 provides more details on the related works.). In practice, however, the generalisation performance of current state-of-the-art CNN methods is mainly limited due to three reasons. First, the EEG presents inherent challenges which make the application of deep learning methods to real-world EEG datasets more difficult. For example, the recorded EEG signals are high-dimensional, non-linear, nonstationary (Cole & Voytek, 2019; Gramfort et al., 2013; Jia et al., 2021), have a low signal-to-noise ratio (Bigdely-Shamlo et al., 2015; Jas et al., 2017), are strongly influenced by artefacts caused by external environmental factors (e.g. electrical interference from external sources) (Islam et al., 2016; Kane et al., 2017) or physiological sources (e.g. cardiac, muscle activity, eye movements, or fatigue) (Britton et al., 2016), and variations in recording protocols and labelling standards within clinical data can occur (Poziomska et al., 2024). In addition, pathological EEG patterns exhibit a wide range of physiological variability across both patients (Nahmias et al., 2019) and neurological conditions (Emmady & Anilkumar, 2023; Nayak & Anilkumar, 2020; Smith, 2005). Due to these challenges, the classification performance can vary between patients (Altuwaijri et al., 2022; Lashgari et al., 2020; Roy et al., 2019b; Schirrmeyer et al., 2017b). Second, due to the scarcity of publicly available datasets, current approaches are mainly trained and evaluated on a single dataset, the TUH Abnormal EEG Corpus (TUAB) (López de Diego, 2017). While it is therefore well established that these methods perform well on this small, homogeneous dataset, this limitation restricts the generalisability of these approaches. For example, recent studies have observed significant differences in the performance of several previously established models after training or validating them on other datasets (Darvishi-Bayazi et al., 2023; Khan et al., 2022; Kiessner et al., 2023; 2024; Nahmias & Kontson, 2020; Poziomska et al., 2024; Van Leeuwen et al., 2019; Western et al., 2021). Lastly, single-scale convolutions, with manually and arbitrarily defined kernel sizes (Emsawas et al., 2022), are mainly used for designing CNNs; however, they are less capable of handling individual differences in the EEG signal. For example, several studies have found that, due to variability in the EEG signal, the optimal kernel size differs from subject to subject and from time to time for the same subject (Altaheri et al., 2023b; Altuwaijri et al., 2022; Jia et al., 2021). While multi-branch, multi-scale<sup>1</sup>, or parallel architectures (Altuwaijri et al., 2022; Belwafi et al., 2017; Ingolfsson et al., 2020; Jia et al., 2021; Riyad et al., 2020; Szegedy et al., 2015; Zhang et al., 2023a) have shown promising results in addressing the challenge of heterogeneity in various EEG classification tasks (Cai et al., 2024; Jia et al., 2021; Siddiqua et al., 2024; Yan et al., 2025; Zhu et al., 2023), only a few attempts have been made to address the challenges of general EEG pathology classification by using CNNs with a set of three convolution scales (Brenner et al., 2024; Roy et al., 2019a; Wu et al., 2021). However, these CNN approaches have achieved classification performance similar to that of other CNNs, with accuracies ranging from 85.10% to 87.10%; they also inherit the disadvantages found in their single-scale counterparts. For example, these results were based on the TUAB, and no attempts have been made to evaluate the performance of these CNNs on a larger, more heterogeneous dataset. In addition, smaller kernel sizes are preferred due to lower computational costs (Emsawas et al., 2022), which, however, tend to learn shorter temporal patterns from faster frequency bands (Cohen, 2014; Jia et al., 2021). This also limits the capability of these CNNs to capture the heterogeneity of the EEG signal (Altaheri et al., 2023b; Emsawas et al., 2022). Therefore, the development of an accurate and reliable CNN solution for general EEG pathology classification requires a CNN network design that

<sup>1</sup>In this work, we will use "multi-scale" and "multi-kernel" interchangeably.

can address the issues of high intra- and inter-subject variability in the EEG signal, data diversity, and the heterogeneity of pathological EEG patterns.

In this work, we propose a novel, multi-branch, multi-scale CNN, called Multi-Branch Multi-Kernel Network (Multi-BK-Net) for general EEG pathology classification that extracts long-term and short-term spatiotemporal EEG features by incorporating five parallel branches within the first convolution-pooling block (see Section 2.1). Each branch employs a temporal convolution layer with a different kernel size that were defined considering five clinically relevant frequency bands, namely delta (1–3 Hz), theta (4–7 Hz), alpha (8–13 Hz), beta (14–30 Hz) and low gamma (30–80 Hz) (Brenner et al., 2024). Aiming at the problem of a small amount of training data and to increase data diversity, we combined the predefined training sets of two publicly available datasets for general EEG pathology classification, the TUH Abnormal EEG Corpus (López de Diego, 2017) and TUH Abnormal Expansion Balanced EEG Corpus (TUABEXB) (Kiessner et al., 2023), to optimise and train our Multi-BK-Net. The hyperparameters of the Multi-BK-Net were optimised using multivariate tree-structured Parzen estimators (TPE) (Bergstra et al., 2011; 2013) from the Optuna library (Akiba et al., 2019b) with respect to the mean validation accuracy and mean validation sensitivity values from a 5-fold cross-validation. Our study presents the following contributions to research on deep end-to-end CNNs for general EEG pathology classification:

- We propose the Multi-BK-Net, a multi-branch, multi-scale CNN designed for general EEG pathology classification. We combine the concepts of multi-scale and multi-branch CNNs with the idea of adapting the temporal kernel size based on five clinically relevant frequency bands to accommodate the challenge of heterogeneity inherent in the EEG signal.
- We evaluate our proposed Multi-BK-Net on the predefined test sets of two publicly available datasets, the TUAB and the TUABEXB datasets. Our Multi-BK-Net achieves mean accuracies of 87.75% and 87.01% and mean sensitivities of 83.10% and 84.25%, respectively, and thus outperforms five architectures used as a baseline (Section 3.4), as well as previously reported state-of-the-art deep learning approaches (Section 3.5) on these datasets, setting a new benchmark.
- Our ablation experiments show that using both multiple temporal kernel lengths and multiple parallel branches in the first block significantly improves the performance of the model, particularly in terms of mean accuracy and mean sensitivity (Section 3.6).
- To interpret the trained Multi-BK-Net, we visualise the learned features using Uniform Manifold Approximation and Projection (UMAP) (McInnes et al., 2018) method (Section 3.7) and observe that the Multi-BK-Net forms distinct and compact clusters of samples from the pathological and non-pathological classes, which partly explains the good performance of this architecture. Additionally, we perform an amplitude gradient analysis (Section 3.8) showing that the model’s pathological prediction is sensitive to localised patterns of amplitude changes in different frequency bands. These patterns align with the current neurophysiological knowledge for pathological EEG patterns (Amin et al., 2023; Emmady & Anilkumar, 2023; Hoppe, 2018; Kane et al., 2017; Medithe & Nelakuditi, 2016; Tatum & William, 2021) and with the pathological patterns that have been identified by human experts in the corresponding clinical EEG reports.
- Overall, our results demonstrate the effectiveness of our multi-branch, multi-scale CNN solution in accurately and reliably classifying general EEG pathology. Our findings contribute to the broader discussion on the challenges of deep learning-based general EEG pathology classification, providing new insights for the design of future end-to-end CNNs. This will ultimately advance the development of more reliable and robust methods for automated classification of general EEG pathology.

## 2 Methods

### 2.1 Multi-Branch Multi-Kernel Network (Multi-BK-Net)

We propose a novel deep, multi-scale, and multi-branch CNN, called Multi-BK-Net, for general EEG pathology classification. This model extracts long-term and short-term spatiotemporal EEG features by incorporating a dedicated multi-branch, multi-scale first block that processes EEG inputs for different temporal

patterns. The block is followed by three standard convolution-mean-pooling blocks and a softmax classification layer. Figure 1 shows the overall structure of the proposed Multi-BK-Net architecture. The first convolutional block is divided into five branches (multi-branches). Each convolutional branch contains two convolutional layers and a mean-pooling layer. The first convolutional layer performs temporal filtering with a set of seven convolutional filters. Each branch uses a different kernel size for temporal filtering (multi-scale). This allows the model to extract and combine features from different time scales. The choice of kernel size was based on five frequency bands, namely, delta (0.5 to 4 Hz), theta (4 to 7 Hz), alpha (8 to 12 Hz), beta (13 to 30 Hz), and low gamma (30 to 80 Hz) (Brenner et al., 2024; Nayak & Anilkumar, 2020), by also considering the sampling rate. In this study, the EEG signal is represented as a 2D input, with the dimension ( $E \times T$ ) and with a window size of 60 seconds and a sampling frequency of 100 Hz, corresponding to 6000 time points, an input shape of  $21 \times 6000$  (for more details, see Section 3.1). Accordingly, we used five different temporal kernel sizes:  $1 \times 200$ ,  $1 \times 25$ ,  $1 \times 13$ ,  $1 \times 7$  and  $1 \times 3$ . In the second layer, a spatial convolution is applied. Each filter performs a spatial convolution with weights for all possible electrode pairs, using the filters from the preceding temporal convolution. Note that there is no activation function between the two layers; therefore, in principle, they could be combined into a single layer. However, we chose a split convolution because it has been shown to outperform a combined temporal-spatial convolutional layer (Schirrneister et al., 2017a;b). After the spatial convolution, mean pooling is applied. As successfully done by (Schirrneister et al., 2017a), the pooling strides were moved directly to the convolutional layers preceding each pooling. Finally, the outputs of each branch are concatenated along the feature axis, resulting in a feature vector of shape (1, 35, 397, 1). The concatenated tensors are then passed through three convolution-mean-pooling blocks. Each convolution-mean-pooling block consists of a convolutional layer and a mean pooling layer. The output of the fourth block is then passed through a softmax classification layer to receive the final prediction. In total, the Multi-BK-Net contains 1,038,683 trainable parameters. We implemented our model in Braindecode (BD), an open-source Python toolbox for decoding raw electrophysiological brain data with deep learning models (Schirrneister et al., 2017b). The Xavier initialisation method (Glorot & Bengio, 2010) was used to initialise the weights. The biases have been initialised to 0. We applied group normalisation (Wu & He, 2018) to the output of the convolutional layers before the nonlinearity. We used Gaussian Error Linear Units, or GELUs ( $GELU(x) = xP(X \leq x) = x\Phi(x)$ ) (Hendrycks & Gimpel, 2023) as activation functions. Additionally, we used dropout as a regularisation strategy (Srivastava et al., 2014). Appendix A.3 provides a more detailed description of the design choices and hyperparameter optimisation procedure.

## 2.2 Network Training

For optimisation, we used the adaptive moment estimation with decoupled weight decay (AdamW) optimiser (Loshchilov & Hutter, 2019) with an optimised learning rate and beta1 parameter and the negative log likelihood loss (Terven et al., 2024). We used cosine annealing (Loshchilov & Hutter, 2016) to schedule the learning rates for both the gradient and weight decay updates, and refrained from restarting the learning rate. We trained the model on non-overlapping, equally sized time windows of size 60 seconds using the trial-wise training strategy<sup>2</sup> as described by Schirrneister et al. (2017b). We repeated the training ten times, using a fixed set of different random seeds to account for model variance, due to random weight initialisation. The final predictions are obtained by first averaging the predictions of each window for each recording, and then calculating the overall performance across recordings (see Section 3.3 for more details). The list of hyperparameters for training is provided in Table 1. The code of this study is available at [anonymous.4open.science/r/Multi-BK-Net-general-EEG-pathology-classification-4FE5](https://anonymous.4open.science/r/Multi-BK-Net-general-EEG-pathology-classification-4FE5).

## 3 Experiments and Results

To demonstrate the effectiveness of our method, we evaluated the Multi-BK-Net on two public datasets for general EEG pathology classification: I) the small, homogeneous, public dataset TUH Abnormal EEG Corpus (TUAB) (López de Diego, 2017) for comparison with previously published accuracies and II) the

<sup>2</sup>During trial-wise training, a complete window is pushed through the network. The network then produces a prediction, which is compared to the target (label) for that window (trial) to compute the loss.

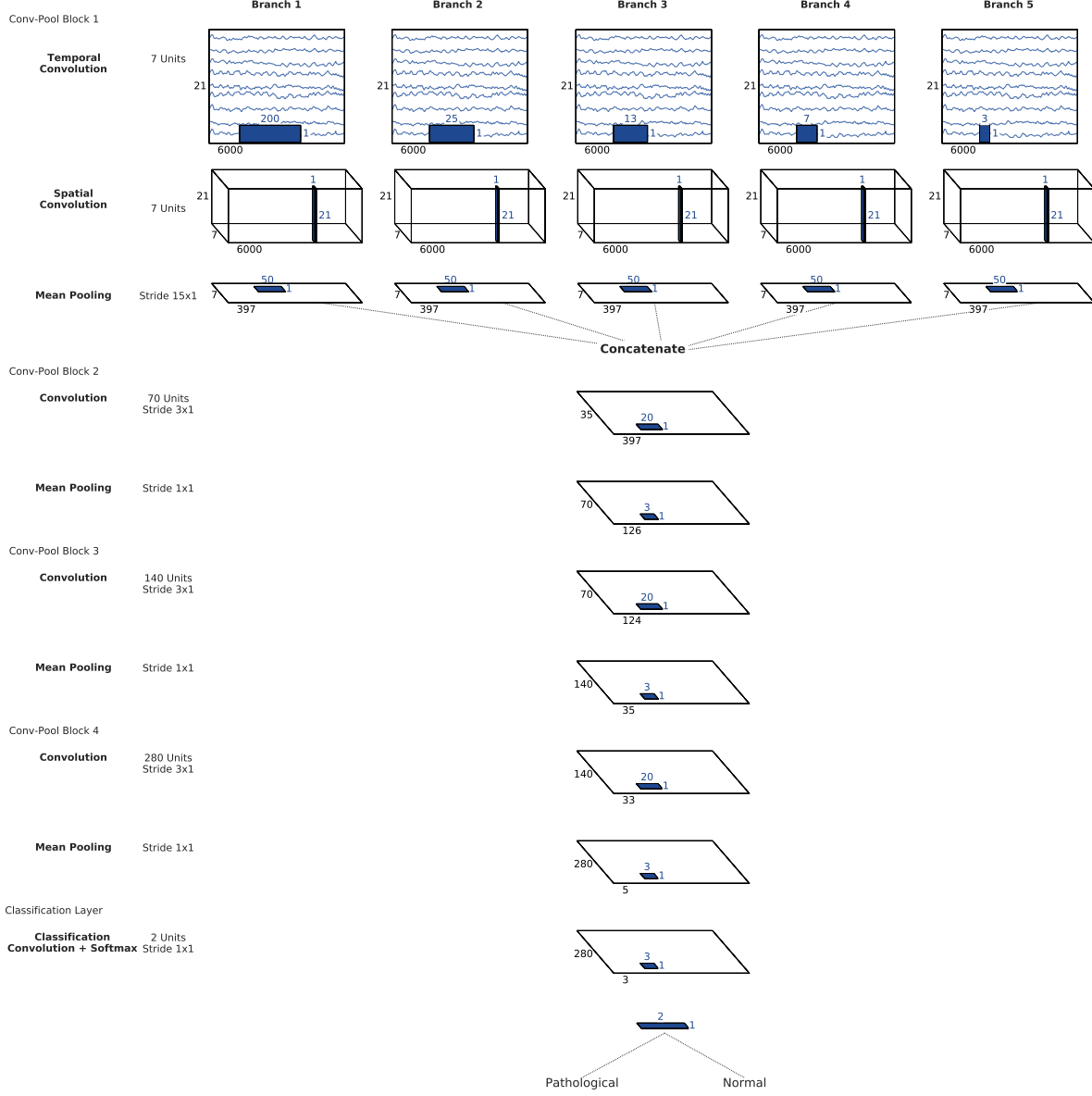


Figure 1: Multi-BK-Net architecture. The first block is divided into five branches, each applying a temporal convolution with a different kernel size, followed by a spatial convolution. There is no activation function between the two convolutional layers; after that, a mean-pooling layer is added. The EEG input (at the top) is passed through each branch. The outputs of all branches are concatenated and then used as input for the subsequent convolution-mean-pooling blocks. Black cuboids: input/feature maps; coloured cuboids: convolution/pooling kernels. The corresponding sizes are shown in black or in colour. Note that the proportions of maps and kernels in this schematic are only approximate.

larger, recently introduced TUH Abnormal Expansion Balanced EEG Corpus (TUABEXB) (Kiessner et al., 2023) to evaluate the classification methods with a larger heterogeneous number of EEG recordings. First, we compared the performance of the Multi-BK-Net to five baseline architectures. Following this, we conducted a performance comparison with previously reported state-of-the-art deep learning approaches. Then, we performed ablation experiments to demonstrate the superiority of the proposed Multi-BK-Net and to further

Table 1: Hyperparameters of the Multi-BK-Net architecture. Hyperparameters have been optimised in preliminary experiments; more details on the hyperparameter optimisation and design choices are provided in Appendix A.3.

Hyperparameter	Selected value
Total number of temporal convolution filters	35 (7 filters per branch)
Normalisation	GroupNorm
Activation functions	GELU
Pooling mode first block	Mean
Pooling mode remaining blocks	Mean
Forth conv-pooling-block	True
Forth conv-pooling-block broader	True
Pool length	3
Pool stride	3
Stride before pool	True
Dropout	0.502959339666169
Filter length convolution blocks	20
Input window size	6000
Weighted loss factor pathological	1
Optimizer	AdamW
Optimizer beta1	0.5
Learning rate	0.0031414364096615
Weight decay	1.8397405899531204e-05
Batch size	64
Number of epochs	42
Number of channels	21

highlight the importance of using multiple temporal kernel sizes and multiple branches in the first block. In addition, we visualised the learned features of the Multi-BK-Net to highlight the robustness of the Multi-BK-Net and provide a comprehensive understanding of model performance. Lastly, we performed an amplitude gradient study to determine how sensitive the Multi-BK-Net’s prediction is to amplitude changes in different frequency bands.

### 3.1 Datasets and Preprocessing steps

The TUAB consists of 2,993 recordings (49.18% pathological) from 2,329 patients (52.09% female, mean age:  $48.55 \pm 17.86$  years) that are divided into a predefined training set (2,717 recordings) and an evaluation set (276 recordings). In contrast, the TUABEXB contains 8,879 recordings (49.75% pathological) obtained from 7,006 patients (mean age:  $47.7 \pm 21.2$  years; 51.7% female) and is divided into a predefined training set (7990 recordings) and an evaluation set (889 recordings). For training, we used the concatenation of the TUAB and TUABEXB training sets, which we refer to as the TUH Abnormal Combined EEG Corpus (TUABCOMB). The TUABCOMB contains 10,707 recordings from 8,549 patients (mean age:  $47.91 \pm 20.53$  years; 52.5% female), of which 5,321 recordings (49.70%) are classified as pathological. We evaluated our models on the predefined test set of both datasets, TUAB and TUABEXB, respectively (see Section 3.3). We applied the following minimal preprocessing steps to the raw EEG data: a) selected a set of 21 electrode positions<sup>3</sup> present in all recordings; b) discarded the first 60 seconds of each recording as they contain stronger artefacts; c) used up to 20 minutes of the remaining recording to speed up the computations; d) clipped the EEG recordings at  $\pm 800 \mu V$  to reduce the effect of strong artefacts; e) re-referenced all recordings to the Common Average Reference (CAR) f) resampled the data to 100 Hz to account for the different sampling rates and to further speed up the computation. To apply CNNs to EEG classification, the raw EEG signals

<sup>3</sup>This set of EEG channels included the channels A1, A2, C3, C4, Cz, F3, F4, F7, F8, Fp1, Fp2, Fz, O1, O2, P3, P4, Pz, T3, T4, T5 and T6.

measured in microvolts ( $\mu\text{V}$ ) can be represented as input as a 2D array with the number of time steps ( $T$ , temporal dimension) as the width and the number of electrode channels ( $E$ , spatial dimension) as the height, resulting in an input shape of  $E \times T$ . To create the inputs for the networks, we process sliding windows over the EEG data,  $D \in \mathbb{R}^{T \times E}$ , where  $T$  represents the recording duration and  $E$  denotes the number of electrodes. Windows are created as input features  $x_t \in \mathbb{R}^I$  with input dimensionality  $I = E \cdot \lfloor T/S \rfloor \cdot f_s$  at time point  $t$ , where  $S$  is the stride and  $f_s$  is the sampling frequency. We use non-overlapping windows of 60 seconds, which corresponds to a stride of  $S = 60$  seconds. The corresponding labels are  $y \in \{0, 1\}^{\lfloor T/S \rfloor}$ , where  $y_i = 0$  indicates a non-pathological window and  $y_i = 1$  indicates a pathological window. A tensor of the shape (batch\_size,  $E$ ,  $T$ ) is given as input to the architectures.

### 3.2 Baselines

To illustrate the superiority of the proposed Multi-BK-Net, we compare it to several approaches reported in previous work. For general EEG pathology classification, it is challenging to compare approaches without evaluating them within the same framework, due to factors such as different evaluation methods, preprocessing steps, training strategies, and several datasets or dataset versions. Therefore, to evaluate the performance of our architecture, we compared our method with the following five baseline architectures:

- Deep4Net (Schirrmester et al., 2017b;a), a CNN consisting of four convolution-max-pooling blocks.
- ShallowNet (Schirrmester et al., 2017b;a), a CNN with a convolution-pooling block that was inspired by the Filter Bank Common Spatial Patterns pipeline (Ang et al., 2008; Chin et al., 2009).
- TCN (Bai et al., 2018; Chrabaszcz, 2018; Gemein et al., 2020), a temporal CNN consisting of five residual blocks.
- EEGNet (Lawhern et al., 2018), a compact CNN with depthwise and separable convolutions.
- ChronoNet (Roy et al., 2019a), a deep recurrent neural network combining multiple Conv1D layers with multiple filters of varying sizes and stacked Gated Recurrent Unit (GRU) layers.

These architectures have been implemented in Braindecode and have also shown high performance on the TUAB and/or TUABEXB in various studies using different preprocessing steps. Additional details on the baseline architectures, the list of hyperparameters for training and details on the training procedure are provided in the Appendix A.4.

### 3.3 Evaluation of Classification Performance

For the final evaluation, we trained the model on the full TUABCOMB training set and evaluated the model on the withheld final evaluation sets, as predefined in the TUAB and TUABEXB, respectively. To manage the statistical variance caused by initialisation and to improve the comparison between training and model configurations, we repeated the training and evaluation ten times for each model and reported the average of evaluation metrics across multiple runs (Bouthillier et al., 2021; Picard, 2023; Wightman et al., 2021). To ensure comparability with previous work, we evaluated the overall performance of a model using the prediction for each recording. The model’s prediction for a recording was calculated by averaging the outputs from all windows of that recording. Each recording was then classified as non-pathological or pathological based on its mean window probability. The performance of the different architectures was evaluated using the following general classification metrics: *accuracy*, *balanced accuracy*, *sensitivity*, *specificity* and *F2-score*. In addition, we used the Mann-Whitney U test (Mann & Whitney, 1947) to test for the statistical significance of the difference in performance metrics between the proposed Multi-BK-Net and each of the baseline architectures (H1). The null hypothesis (H0: No significant differences in the performance of Multi-BK-Net compared to a baseline architecture) was rejected with a p-value of  $p < 0.05$ . To correct for multiple testing, we additionally performed a Bonferroni correction (Bonferroni, 1936) with  $\alpha = 0.05$  for all performance comparisons involving our Multi-BK-Net and each baseline architecture on the corresponding evaluation set.

### 3.4 Performance Comparison to Baseline Architectures

In this section, we examine the performance of the proposed Multi-BK-Net in comparison to five baseline architectures on two predefined test datasets, the TUAB and the TUABEXB. The classification performance of the Multi-BK-Net and all baseline architectures on the TUAB is shown in Figure 2. With a mean accuracy of 87.75%, the Multi-BK-Net shows the best classification performance. This was also statistically significantly better than the baseline models ( $p < 0.001$ , Mann-Whitney U test). However, in medical diagnosis, such as EEG pathology classification, it is more important to identify all potential pathological cases (high recall/sensitivity) even at the expense of some false positives (lower precision). In this context, the F2-score is a suitable additional metric to compare the different architectures. With a mean sensitivity of 83.10% and a mean F2-score of 84.05%, Multi-BK-Net statistically significantly outperformed all other architectures ( $p < 0.001$ , Mann-Whitney U test), indicating better performance in terms of recall. Figure 3 compares the classification performance of the Multi-BK-Net and five baseline architectures on the TUABEXB evaluation set. Again, the Multi-BK-Net outperformed the baseline models. In particular, Multi-BK-Net achieved a higher mean accuracy of 87.01% on the dedicated test set, which was also statistically significantly different from the baseline CNNs ( $p < 0.001$ , Mann-Whitney U test). We can also see that the Multi-BK-Net achieved statistically significantly better mean sensitivity (84.25%) and mean F2-score (85.17%) than each of the baseline CNNs ( $p < 0.001$ , Mann-Whitney U test). Additional classification results are presented in Appendix A.6.1. In general, our method achieved higher classification performance than five baseline architectures and was more effective in identifying pathological recordings. This suggests that Multi-BK-Net could be used as an alternative to CNN for classifying general EEG pathology.

### 3.5 Performance Comparison to Previously Reported State-of-the-art Approaches

To further evaluate the Multi-BK-Net, we compared its classification performance with that of other state-of-the-art end-to-end approaches reported in previous studies. More details on the previously published studies are provided in Appendix A.6.2. Table 2 and Table 3 summarise the comparison between our Multi-BK-Net and other methods on the TUAB and TUABEXB datasets, respectively. Overall, the Multi-BK-Net achieved a similar or higher level of performance than other deep learning approaches in all three evaluation metrics, with small advantages for the Multi-BK-Net in some settings. In particular, for the TUAB evaluation set, the following observations can be made by further analysis of the Table 2: First, compared to the four single-scale CNNs, namely Deep4Net, ShallowNet, TCN and EEGNet (Darvishi-Bayazi et al., 2023; Gemein et al., 2020; Khan et al., 2022; Kiessner et al., 2023; 2024; Schirrmeister et al., 2017a; Western et al., 2021), Multi-BK-Net improves the mean classification accuracy by 1.85%, 2.24%, 1.26% and 3.08%, respectively. Second, compared to the three architectures using a set of three convolutional scales, i.e. ChronoNet (Roy et al., 2019a), XceptionTime model (Brenner et al., 2024) and IRCNN (Wu et al., 2021), Multi-BK-Net improves the classification accuracy by 1.18%, 1.65% and 0.65%, respectively. In other words, the Multi-BK-Net achieved a classification accuracy of 87.75%, thus outperforming the best previously reported state-of-the-art performance of the IRCNN (Wu et al., 2021). Third, Multi-BK-Net also demonstrates higher classification accuracy than two pre-trained transformer-based foundation models, BIOT (3.2M parameters) and LaBraM (396M parameters), both of which require a substantial amount of additional EEG data for pre-training and substantial computational resources due to their large number of parameters. Finally, perhaps the most clinically important finding is that, with a mean sensitivity of 83.10%, Multi-BK-Net outperforms all other approaches, reporting mean sensitivity averaged over multiple runs, by at least 3.26% (see Table 2). Furthermore, for the TUABEXB evaluation set, we have the following observations from further analysis of Table 3: Compared to previous results (Kiessner et al., 2023; 2024), the Multi-BK-Net achieved a better mean accuracy of 87.01% and a mean sensitivity of 84.25%, representing improvements of at least 1.66% and 2.76%, respectively. Thus, Multi-BK-Net achieved new state-of-the-art results on the predefined TUABEXB evaluation set, setting a new benchmark. Additional results and a more detailed discussion are provided in Appendix A.6.2. Overall, we can see that the Multi-BK-Net outperforms previously reported results on both datasets, in terms of mean accuracy. In addition, Multi-BK-Net achieved a higher mean sensitivity of 83–84%, indicating improved performance for real-world medical applications, particularly for medical screening methods where high sensitivity is crucial for accurately identifying the presence of EEG pathology.



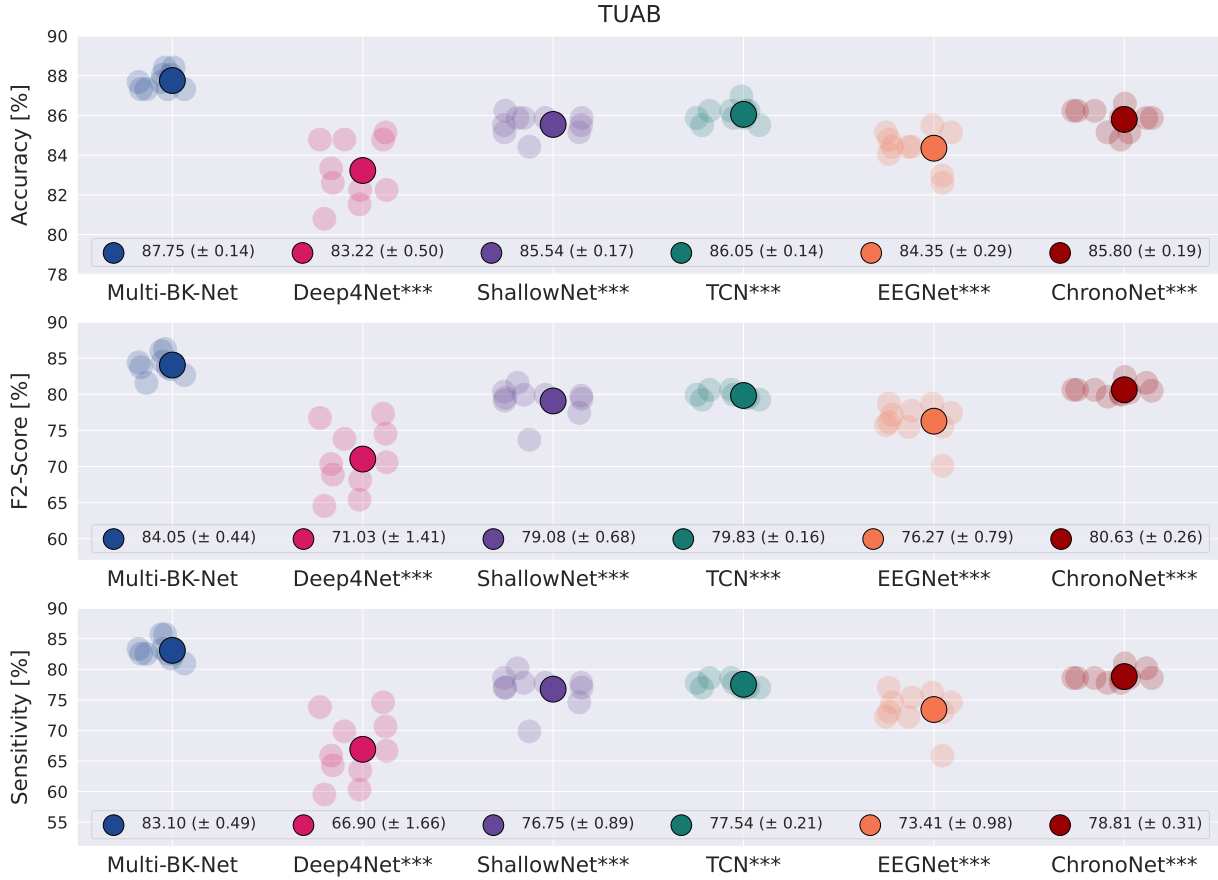


Figure 2: Performance comparison between our proposed Multi-BK-Net and five baseline architectures on the predefined TUAB evaluation set. Each transparent marker represents the performance of a single run, and each larger, bold symbol represents the mean performance score averaged across ten independent runs. The mean standard error is given in parentheses. Stars indicate statistically significant differences in performance score between the corresponding baseline architecture and the Multi-BK-Net (Bonferroni-corrected two-sided Mann-Whitney U test,  $p < 0.05$ : \*,  $p < 0.01$ : \*\*,  $p < 0.001$ : \*\*\*). For additional results, see Figure 6, Appendix A.6.1.

### 3.6 Ablation Experiments

To further validate the effectiveness of the proposed Multi-BK-Net and to highlight the importance of the multi-scale and multi-branch components in the Multi-BK-Net, we performed several ablation experiments, where we a) compared the performance of the Multi-BK-Net to larger variants of the Deep4Net having the same width or approximately the same size as the Multi-BK-Net ("larger Deep4Net variants"); b) removed four of the five branches of the Multi-BK-Net (Single-Temporal-Kernel Single-Branch Net, STKSBN, "single temporal kernel + single branch") while changing the number of filters from 7 to 35; and c) used five different temporal kernel lengths but one branch, i.e. concatenating the output of the five temporal convolution layers before the spatial convolution layer (Multi-Temporal-Kernel Single-Branch Net, MTKSBN, "multiple temporal kernels + single branch"). The results of these experiments can be seen in Table 4 and Table 5, while Appendix A.6.4 provides more details on the ablation experiments, the size of the models (see Table 14) and a more detailed discussion of the results. In comparison to the Multi-BK-Net, "larger Deep4Net variants" achieved statistically significantly lower mean accuracies and mean sensitivities (Table 4), indicating that the superiority of the Multi-BK-Net over the baseline models was not due to the increased model size alone. When compared against our Multi-BK-Net (Table 5), the networks with "single

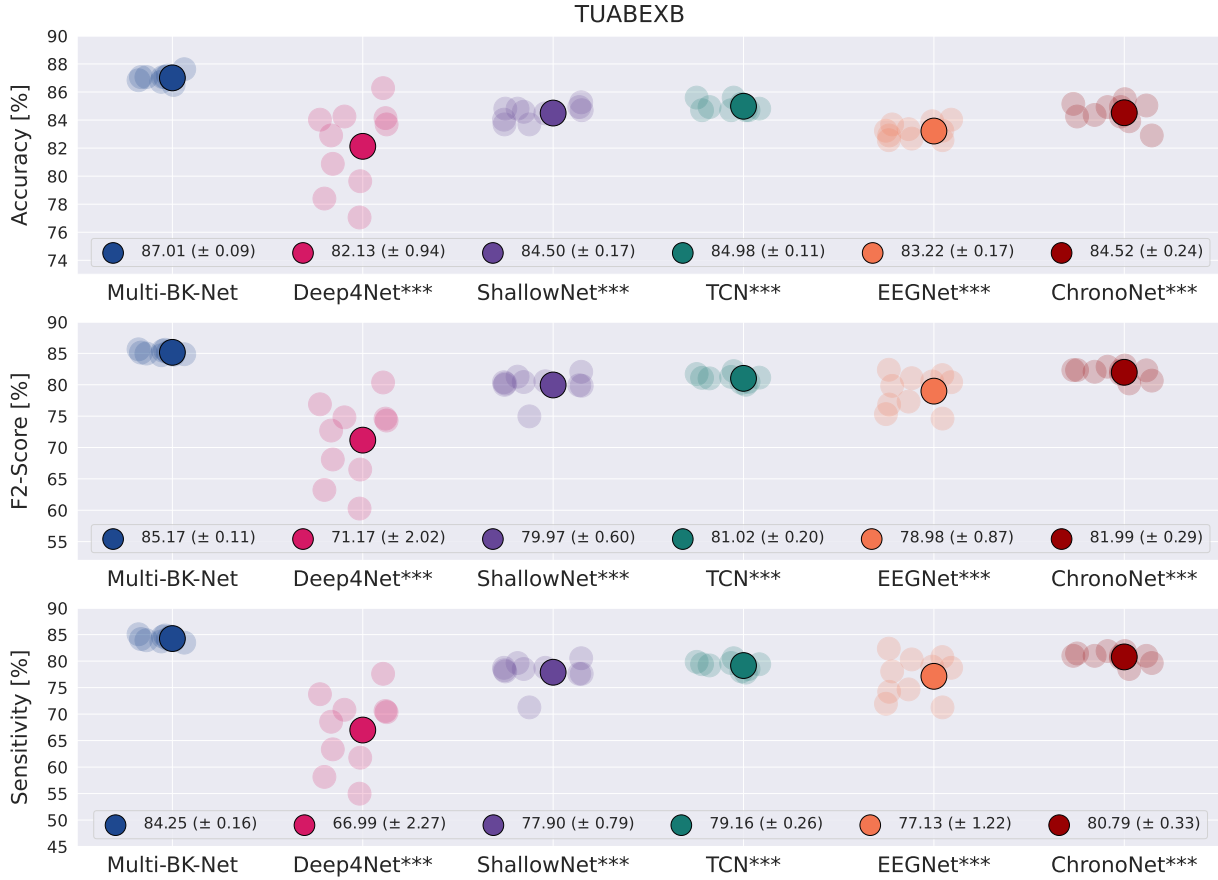


Figure 3: Performance comparison between our proposed Multi-BK-Net and five baseline architectures on the predefined TUABEXB evaluation set. Conventions as in Figure 2. For additional results, see Figure 7, Appendix A.6.1.

temporal kernel + single branch" and "multiple temporal kernel + single branch" suffered from a statistically significant drop in classification performance, especially in terms of mean accuracies and mean sensitivities, thereby highlighting the importance of these components towards classification.

### 3.7 Network Interpretation using a UMAP Visualisation

To interpret the proposed network, we extracted the learned features of the last convolution-pooling block of the Multi-BK-Net and visualised them using the Uniform Manifold Approximation and Projection (UMAP) (McInnes et al., 2018) method, which is a dimensionality reduction technique well-suited for visualising high-dimensional data, such as feature representation vectors. Figure 4 shows the UMAP visualisation of the learned features. As can be seen in the plot, the representations for the pathological and non-pathological samples form distinct and compact clusters with minimal overlap, indicating that our proposed Multi-BK-Net is very robust in classifying pathological samples and extracts well-classifiable features, which partly explains the good performance of this architecture. Further analysis in Appendix A.6.3 reveals that compared to Multi-BK-Net, the three baseline models (Deep4Net, TCN and ChronoNet) tend to form smaller and more dispersed clusters, possibly indicating greater within-class variability or a less precise representation of the pathological class. These results further align with the quantitative performance metrics, such as accuracy and sensitivity, shown in Section 3.4. Overall, the results of the UMAP visualisation of the learned features further support the effectiveness of our method in extracting features and identifying pathological samples.

Table 2: Performance comparisons of the Multi-BK-Net with previously reported state-of-the-art deep learning methods on the TUAB dataset. Accuracy (Acc.), sensitivity (Sens.) and specificity (Spec.) scores are given in %, n.a.: not available. Results marked with  $\bowtie$  are reported as balanced accuracy.  $R$ : Publicly available reimplementation. Studies marked with  $\times$  used data augmentation. Studies marked with  $\triangleleft$  used only the first minute of each recording for performance evaluation. Studies marked with  $\dagger$  used a cropped training strategy as described by Schirrmeister et al. (2017a;b). Studies marked with  $\emptyset$  reported mean results averaged across several independent runs. Results marked with  $\dagger$  are the performance on TUAB using the TUABEXB training (Kiessner et al., 2023) instead of the TUAB training set for training. Results marked with  $\ddagger$  are the performance on TUAB using the TUABCOMB (Kiessner et al., 2024) for training.

Study	Architecture	ACC [%]	SENS[%]	SPEC [%]
Schirrmeister et al. (2017a) $\dagger \emptyset$	Deep4Net	85.40	75.12	94.13
	ShallowNet	84.50	77.32	90.53
Roy et al. (2018) $\triangleleft \emptyset$	1D-CNN-RNN	82.27	n.a.	n.a.
Roy et al. (2019a) $\emptyset$	ChronoNet	86.57	n.a.	n.a.
Gemein et al. (2020) $\dagger \emptyset$	TCN	86.10	79.70	91.60
	Deep4Net	84.60	75.90	91.90
	ShallowNet	84.10	79.70	87.90
	EEGNet	83.40	72.10	92.90
Western et al. (2021) $\emptyset$	Deep4Net	85.90	77.00	93.30
Wu et al. (2021) $\triangleleft \emptyset$	IRCNN	87.10	n.a.	n.a.
Khan et al. (2022)	ChronoNet $R$	81.00	n.a.	n.a.
	ShallowNet	85.00	n.a.	n.a.
	Deep4Net	84.00	n.a.	n.a.
	Hybrid CNN+LSTM	85.00	n.a.	n.a.
Darvishi-Bayazi et al. (2023) $\times \emptyset$	Deep4Net	81.64 $\bowtie$	n.a.	n.a.
	ShallowNet	82.40 $\bowtie$	n.a.	n.a.
	TCN	81.69 $\bowtie$	n.a.	n.a.
	EEGNet	81.40 $\bowtie$	n.a.	n.a.
Kiessner et al. (2023) $\dagger \emptyset$	Deep4Net	85.51	75.95	93.53
	ShallowNet	84.13	79.84	87.73
	TCN	85.72	78.81	91.53
	EEGNet	84.67	72.94	94.53
	Deep4Net $\dagger$	84.89	70.95	96.60
	ShallowNet $\dagger$	85.51	77.30	92.40
	TCN $\dagger$	86.49	77.06	94.40
	EEGNet $\dagger$	83.51	77.22	88.80
Yang et al. (2023) $\bowtie \emptyset$	BIOT	79.59 $\bowtie$	n.a.	n.a.
Kiessner et al. (2024) $\dagger \emptyset$	Deep4Net $\ddagger$	85.29	73.73	95.00
	ShallowNet $\ddagger$	85.22	76.83	92.27
	TCN $\ddagger$	86.09	77.70	93.13
	EEGNet $\ddagger$	83.80	78.41	88.33
Brenner et al. (2024) $\triangleleft$	XceptionTime model	85.10	85.70	84.70
Jiang et al. (2024) $\bowtie \emptyset$	LaBraM-Huge	82.58 $\bowtie$	n.a.	n.a.
<b>Proposed method <math>\emptyset</math></b>	<b>Multi-BK-Net</b>	<b>87.75</b>	<b>83.10</b>	<b>91.93</b>

### 3.8 Network Interpretation using an Amplitude Gradient Analysis

In this section, we investigate how sensitive the Multi-BK-Net’s pathological prediction is to amplitude changes in the input. Previous neurophysiological studies and glossaries have described the most common abnormal EEG patterns, for example, amplitude changes within multiple frequency bands (Amin et al., 2023; Emmady & Anilkumar, 2023; Hoppe, 2018; Kane et al., 2017; Medithe & Nelakuditi, 2016; Tatum & William,

Table 3: Performance comparisons of the Multi-BK-Net with previously reported state-of-the-art deep learning methods on the TUABEXB dataset. Accuracy (Acc.), sensitivity (Sens.) and specificity (Spec.) scores are given in %, n.a.: not available. Conventions as in Table 2.

Study	Architecture	ACC [%]	SENS[%]	SPEC [%]
Kiessner et al. (2023) ‡ ∅	Deep4Net	83.94	71.99	95.75
	ShallowNet	84.47	79.32	89.55
	TCN	83.73	77.49	89.91
	EEGNet	82.38	80.79	83.96
Kiessner et al. (2024) ‡ ∅	Deep4Net ‡	85.35	75.27	95.32
	ShallowNet‡	84.58	79.79	89.31
	TCN‡	85.08	79.12	90.98
	EEGNet‡	82.65	81.49	83.80
<b>Proposed method ∅</b>	<b>Multi-BK-Net</b>	<b>87.01</b>	<b>84.25</b>	<b>89.73</b>

Table 4: Ablation experiment "larger Deep4Net variants". Comparison of the classification performance of Multi-BK-Net with larger Deep4Net variants on the predefined TUAB and TUABEXB evaluation sets. Deep4Net as defined in Schirrmeyer et al. (2017a;b). Deep4Net35 (No. start filters: 35); Deep4Net47 (No. start filters: 47); Deep4Net48 (No. start filters: 48). Classification metrics are averaged across ten independent runs. The mean standard error is given in parentheses. ACC: Accuracy, BACC: Balanced accuracy, F2: F2-score, SENS: Sensitivity, SPEC.: Specificity. Stars indicate statistically significant differences in performance between the corresponding model variant and the Multi-BK-Net (Bonferroni-corrected, one-sided Mann-Whitney U test,  $p < 0.05$ : \*,  $p < 0.01$ : \*\*,  $p < 0.001$ : \*\*\*).

Dataset	Architecture	ACC. [%]	BACC. [%]	F2 [%]	SENS. [%]	SPEC. [%]
TUAB	Deep4Net	83.22***	81.92***	71.03***	66.90***	96.93
	Deep4Net35	84.78***	83.63***	74.21***	70.40***	96.87
	Deep4Net47	84.89***	83.90***	75.72***	72.46***	95.33
	Deep4Net48	84.96***	83.94***	75.52***	72.14***	95.73
	<b>Multi-BK-Net</b>	<b>87.75</b>	<b>87.36</b>	<b>84.05</b>	<b>83.10</b>	<b>91.93</b>
TUABEXB	Deep4Net	82.13***	82.04***	71.17***	66.99***	97.09
	Deep4Net35	83.52***	83.45***	74.59***	70.88***	96.02
	Deep4Net47	84.50***	84.44***	76.73***	73.42***	95.46
	Deep4Net48	83.87**	83.80**	75.14***	71.52***	96.09
	<b>Multi-BK-Net</b>	<b>87.01</b>	<b>86.99</b>	<b>85.17</b>	<b>84.25</b>	<b>89.73</b>

2021). Hence, amplitude changes are useful for EEG classification. To determine how sensitive the Multi-BK-Net prediction is to amplitude changes in different frequency bands, we calculated the gradients of the output with respect to the amplitudes of the input, i.e. the recordings of both the TUAB and TUABEXB evaluation sets (Ancona et al., 2018; 2019; Gemein et al., 2020; Schirrmeyer et al., 2017a). The gradients were then grouped by pathological status and averaged over all ten runs. The gradients of the model prediction with respect to the input amplitudes linearly approximate how the prediction of a class behaves in response to changes in the inputs. For example, if the amplitude at a certain channel is changed by a value  $x$ , the model output for the class pathological will change approximately by  $\text{gradient} \cdot x$ . The resulting patterns linearly approximate how changes in input amplitude affect the model’s prediction of pathology, indicating which changes in input amplitudes are most informative or interesting with respect to the task of classifying general EEG pathology. In addition, we analysed the clinical EEG reports, which are included in the TUAB and TUABEXB datasets. Figure 5 shows scalp maps of the amplitude gradients across different frequency bands for the pathological class. The largest absolute values are observed in the delta (0-4 Hz), theta (4-7 Hz) and alpha (8-12 Hz) frequency bands, indicating that the model’s pathological prediction is most sensitive to amplitude changes in these specific frequency bands. In particular, in the delta and theta frequency bands, peaks for positive gradients can be seen at temporal (T3, T4) electrode locations. Specifically, an increase in

Table 5: Ablation experiments "single temporal kernel + single branch" and "multiple temporal kernels + single branch". Comparison of the classification performance of STKSBN variants, MTKSBN variant and Multi-BK-Net on the predefined TUAB and TUABEXB evaluation set. STKSBN uses a single kernel size (e.g., STKSBN7: kernel size of 7) and a single branch in the first block. MTKSBN uses five different kernel sizes, but one branch in the first block. Conventions as in Table 4.

Dataset	Architecture	ACC. [%]	BACC. [%]	F2 [%]	SENS. [%]	SPEC. [%]
TUAB	STKSBN3	84.71***	84.33***	81.00***	79.92***	88.73**
	STKSBN7	84.13***	83.76***	80.46**	79.44**	88.07***
	STKSBN10	84.96***	84.57***	81.20**	80.08*	89.07**
	STKSBN13	86.45**	85.99**	82.17*	80.71**	91.27
	STKSBN25	86.27**	85.81**	82.00**	80.56**	91.07
	STKSBN200	86.38***	85.92***	82.09**	80.63**	91.20
	MTKSBN	86.59**	86.16**	82.55*	81.19*	91.13
	<b>Multi-BK-Net</b>	<b>87.75</b>	<b>87.36</b>	<b>84.05</b>	<b>83.10</b>	<b>91.93</b>
TUABEXB	STKSBN3	82.73***	82.72***	81.48***	80.93***	84.52***
	STKSBN7	84.85***	84.83***	82.60***	81.54***	88.12**
	STKSBN10	86.09**	86.07**	83.76***	82.62***	89.51
	STKSBN13	85.94**	85.82**	83.58**	82.44**	89.40
	STKSBN25	85.95***	85.93***	83.65***	82.53***	89.33
	STKSBN200	86.16**	86.14**	83.62***	82.38***	89.91
	MTKSBN	86.46	86.44	84.26*	83.19*	89.69
	<b>Multi-BK-Net</b>	<b>87.01</b>	<b>86.99</b>	<b>85.17</b>	<b>84.25</b>	<b>89.73</b>

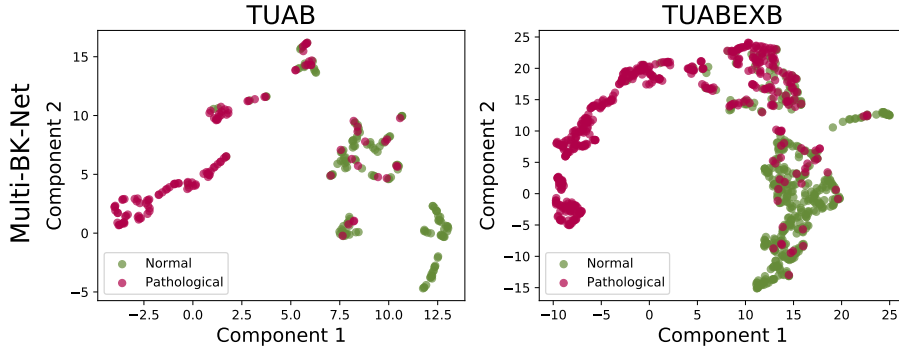


Figure 4: UMAP visualisation of feature representations learned with Multi-BK-Net for pathological and non-pathological recordings on the predefined TUAB (left column) and the TUABEXB (right column) evaluation sets. Pink dots represent the pathological class, while green dots represent the non-pathological class.

amplitude in the temporal areas led to an increase in the prediction of the pathological class. This localised pattern is also consistent with the clinical EEG reports, which often mention that recordings are classified as pathological because of the presence of non-epileptiform abnormalities, such as diffuse or focal slowing in the theta or delta bands (Emmady & Anilkumar, 2023; Nayak & Anilkumar, 2020). Moreover, the clinical EEG reports often cite the "absence of posterior dominant rhythm" as a reason for pathological labelling. Similarly, a negative gradient peak in the alpha band was observed in the occipital brain region, indicating that the predicted probability of the pathological class decreases as the amplitude of the alpha frequency at the occipital electrode locations increases, which, in turn, suggests the presence of the PDR. Overall, the

results align with current neurophysiological knowledge on pathological EEG patterns, pathological patterns identified by human experts in clinical EEG reports, as well as recent studies that implicate the role of low-frequency oscillations in CNNs’ classification of general EEG pathology (Gemein et al., 2020; Nahmias & Kontson, 2020; Schirrmeister et al., 2017a).

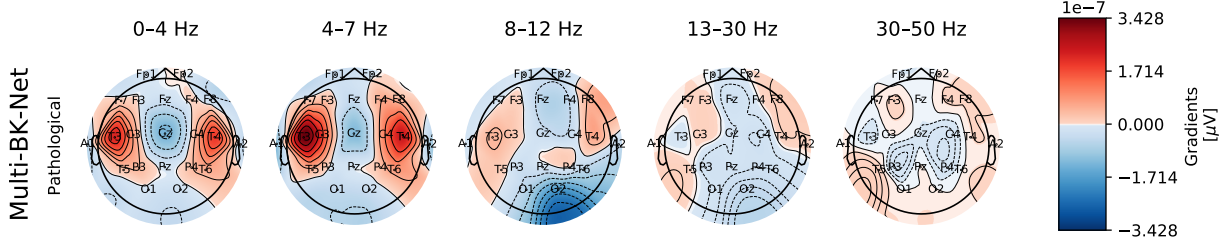


Figure 5: Gradients with respect to the input amplitudes for the pathological class with respect to different frequency bands for the Multi-BK-Net on the combined TUAB and TUABEXB evaluation sets. Delta (0-4 Hz), theta (4-7 Hz), alpha (8-12 Hz), beta (13-30 Hz) and low gamma (30-50 Hz) frequency bands. The red colour indicates positive gradients, while the blue colour indicates negative gradients for the pathological class. Scalp maps are scaled to the absolute maximum value. However, the amplitude gradient results reflect the behaviour of the trained model, and any interpretation of the data itself must be made carefully (Schirrmeister et al., 2017b).

## 4 Discussion and Conclusion

In this work, we demonstrated that through the incorporation of a) multiple temporal kernel lengths based on five clinically relevant frequency bands and b) multiple parallel branches within the first convolution-pooling block, our proposed CNN can outperform five baseline methods in classifying general EEG pathology, while being more effective at capturing the heterogeneity of pathological EEG recordings than their single-scale counterparts. To that effect, our Multi-BK-Net extracted task-discriminative long-term and short-term spatiotemporal EEG features for general EEG pathology classification, achieving higher performance levels than previously reported comparable deep end-to-end state-of-the-art methods on two public datasets. The ablation experiments further highlight the effectiveness of the multi-branch, multi-scale components in our proposed CNN architecture. Additionally, the UMAP visualisation of the learned features shows that the Multi-BK-Net forms more compact and distinct features than the baseline CNNs, which partly explains the good performance of this architecture. Moreover, the correspondence between the model’s sensitivity to localised patterns of amplitude changes in different frequency bands and both current neurophysiological knowledge of pathological EEG patterns and the pathological patterns identified by human experts in the medical report increased the reliability of the method.

This work addresses an important gap in CNN approaches, which has hindered their development into a robust and generalisable approach for deep learning-based EEG pathology classification. While CNN-based methods have shown state-of-the-art performance on the task of general EEG pathology classification (Darvishi-Bayazi et al., 2023; Gemein et al., 2020; Khan et al., 2022; Kiessner et al., 2023; 2024; Van Leeuwen et al., 2019; Western et al., 2021; Wu et al., 2021), their generalisation performance is limited mainly due to the high intra- and inter-subject variability of the EEG signal (Lashgari et al., 2020; Nahmias et al., 2019; Schirrmeister et al., 2017b), the heterogeneity of general EEG pathology patterns (Emmady & Anilkumar, 2023; Nayak & Anilkumar, 2020) and the use of one limited evaluation set (Kiessner et al., 2023; Poziomska et al., 2024). We demonstrated here that an appropriately constructed CNN, incorporating a multi-scale and multi-branch network design, can achieve more accurate and reliable classification performance than state-of-the-art CNN methods. This is evident through the high mean evaluation accuracies and sensitivities achieved by our method on two public datasets, as well as shown in our ablation experiments, thus highlighting its practical effectiveness in classifying general EEG pathology. Our model addresses the challenges inherent in EEG, including the heterogeneity of the EEG signal.

Further practical gains of our method are demonstrated by an improvement in evaluation sensitivity. Our method achieved mean sensitivities of 83.10% and 84.25% on the TUAB and TUABEXB, respectively, outperforming all other approaches by 3.26% and 2.76%, while achieving high classification accuracies. The significance of this result lies in applications such as deep learning-based EEG classification or clinical decision support systems in clinical practice, for which a high and robust sensitivity is crucial to ensure the accurate identification of EEG pathology (Gemein et al., 2020). To be applied in real-world clinical practice, the proposed method must be reliable and utilise robust, well-classifiable features to achieve high performance across datasets. That is why the evaluation on a larger, more heterogeneous dataset was central in this work. Moreover, we interpreted the model using an amplitude gradient study, which shows that the model’s pathological prediction is sensitive to localised patterns of amplitude changes, aligning well with current neurophysiological knowledge of abnormal EEG patterns (Emmady & Anilkumar, 2023; Kane et al., 2017; Nayak & Anilkumar, 2020) and pathological patterns mentioned in clinical EEG reports. In addition, the results of the UMAP visualisation of the learned features further support the effectiveness of our method in extracting features and identifying pathological samples.

This paper describes our efforts to improve the generalisation performance and robustness of CNN methods for general EEG pathology classification. To this end, we proposed the Multi-BK-Net, a novel network architecture with multiple temporal kernels and multiple branches in the first block to suit the challenges of variability and heterogeneity inherent in EEG signals. Further performance improvements could be achieved by exploiting advances in the training of CNNs. Some of these include the incorporation of (self-)attention methods (Altaheri et al., 2023a; Liu et al., 2024; Petit et al., 2021) and increasing the size and diversity of the data by using data from multiple sources (Aerts et al., 2017; Poziomska et al., 2024; Schinkel et al., 2023). As human experts also consider patient-specific information (e.g. age, medication, clinical history) during EEG analysis (Beuchat et al., 2021; Kane et al., 2017; Limotai et al., 2020), future research could extend our methods to accommodate multi-modal inputs, leveraging approaches similar to those explored by Joo et al. (2023), Samak et al. (2023) and Thapa et al. (2024). Moreover, integrating ratings from multiple experts (Stephansen et al., 2018) to evaluate the model performance or using averaged ratings, so-called "fuzzy labels", which represent a probability of pathology, for training might further enhance performance. A systematic optimisation of labelling quality might even help to overcome the current asymptotic limits of predictive accuracy observed in general EEG pathology classification (Darvishi-Bayazi et al., 2023; Kiessner et al., 2024; Poziomska et al., 2024), and hence represent promising avenues for exploration. Overall, our results demonstrate the effectiveness of multi-branch, multi-scale CNN solutions in classifying general EEG pathology by discovering task-relevant features in the recorded brain activity, which can better capture the heterogeneity in pathological samples (Altaheri et al., 2023b; Altuwaijri et al., 2022; Cai et al., 2024; Jia et al., 2021; Liu & Yang, 2021; Liu et al., 2023; Riyad et al., 2020; Siddiqua et al., 2024; Yang et al., 2021; Zhu et al., 2023). By that, our method is a better fit for deep learning-based general EEG pathology classification than traditional CNN methods and thus has significant potential to improve the practicality of deep learning-based general EEG pathology classification.

## References

- Marc Aerts, Girma Minalu, Stefan Bösner, Frank Buntinx, Bernard Burnand, Jörg Haasenritter, Lilli Herzig, J. André Knottnerus, Staffan Nilsson, Walter Renier, Carol Sox, Harold Sox, and Norbert Donner-Banzhoff. Pooled individual patient data from five countries were used to derive a clinical prediction rule for coronary artery disease in primary care. *Journal of Clinical Epidemiology*, 81:120–128, 2017. ISSN 0895-4356. doi: <https://doi.org/10.1016/j.jclinepi.2016.09.011>. URL <https://www.sciencedirect.com/science/article/pii/S089543561630484X>.
- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD ’19, pp. 2623–2631, New York, NY, USA, 2019a. Association for Computing Machinery. ISBN 9781450362016. doi: 10.1145/3292500.3330701. URL <https://doi.org/10.1145/3292500.3330701>.
- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD Interna-*

- tional Conference on Knowledge Discovery and Data Mining*, 2019b.
- Mohammad Ahmad A. Al-Ja'afreh, Geev Mokryani, and Bilal Amjad. An enhanced cnn-lstm based multi-stage framework for pv and load short-term forecasting: Dso scenarios. *Energy Reports*, 10: 1387–1408, 2023. ISSN 2352-4847. doi: <https://doi.org/10.1016/j.egy.2023.08.003>. URL <https://www.sciencedirect.com/science/article/pii/S2352484723011447>.
- Hamdi Altaheri, Ghulam Muhammad, and Mansour Alsulaiman. Physics-informed attention temporal convolutional network for eeg-based motor imagery classification. *IEEE Transactions on Industrial Informatics*, 19(2):2249–2258, 2023a. doi: 10.1109/TII.2022.3197419.
- Hamdi Altaheri, Ghulam Muhammad, Mansour Alsulaiman, Syed Umar Amin, Ghadir Ali Altuwaijri, Wadood Abdul, Mohamed A Bencherif, and Mohammed Faisal. Deep learning techniques for classification of electroencephalogram (eeg) motor imagery (mi) signals: A review. *Neural Computing and Applications*, 35(20):14681–14722, 2023b.
- Ghadir Ali Altuwaijri, Ghulam Muhammad, Hamdi Altaheri, and Mansour Alsulaiman. A multi-branch convolutional neural network with squeeze-and-excitation attention blocks for eeg-based motor imagery signals classification. *Diagnostics*, 12(4), 2022. ISSN 2075-4418. doi: 10.3390/diagnostics12040995. URL <https://www.mdpi.com/2075-4418/12/4/995>.
- Ushtar Amin, Fábio A. Nascimento, Ioannis Karakis, Donald Schomer, and Selim R. Benbadis. Normal variants and artifacts: Importance in eeg interpretation. *Epileptic Disorders*, 25(5):591–648, 2023. doi: <https://doi.org/10.1002/epd2.20040>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/epd2.20040>.
- Ghita Amrani, Amina Adadi, Mohammed Berrada, Zouhayr Souirti, and Saïd Boujraf. Eeg signal analysis using deep learning: A systematic literature review. In *2021 Fifth International Conference On Intelligent Computing in Data Sciences (ICDS)*, pp. 1–8, 2021. doi: 10.1109/ICDS53782.2021.9626707.
- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks, 2018. URL <https://arxiv.org/abs/1711.06104>.
- Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Gradient-based attribution methods. In *Explainable AI: Interpreting, explaining and visualizing deep learning*, pp. 169–191. Springer, 2019.
- Kai Keng Ang, Zheng Yang Chin, Haihong Zhang, and Cuntai Guan. Filter bank common spatial pattern (fbcs) in brain-computer interface. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pp. 2390–2397, 2008. doi: 10.1109/IJCNN.2008.4634130.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- Kais Belwafi, Fakhreddine Ghaffari, Ridha Djemal, and Olivier Romain. A hardware/software prototype of eeg-based bci system for home device control. *Journal of Signal Processing Systems*, 89:263–279, 2017.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL [https://proceedings.neurips.cc/paper\\_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf).
- James Bergstra, Daniel Yamins, and David Cox. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *International conference on machine learning*, pp. 115–123. PMLR, 2013.
- Isabelle Beuchat, Senubia Alloussi, Philipp S Reif, Nora Sterlepper, Felix Rosenow, and Adam Strzelczyk. Prospective evaluation of interrater agreement between eeg technologists and neurophysiologists. *Scientific reports*, 11(1):13406, 2021. doi: <https://doi.org/10.1038/s41598-021-92827-3>.



- Nima Bigdely-Shamlo, Tim Mullen, Christian Kothe, Kyung-Min Su, and Kay A. Robbins. The prep pipeline: standardized preprocessing for large-scale eeg analysis. *Frontiers in Neuroinformatics*, 9, 2015. ISSN 1662-5196. doi: 10.3389/fninf.2015.00016. URL <https://www.frontiersin.org/journals/neuroinformatics/articles/10.3389/fninf.2015.00016>.
- Carlo Bonferroni. Teoria statistica delle classi e calcolo delle probabilita. *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, 8:3–62, 1936.
- Xavier Bouthillier, Pierre Delaunay, Mirko Bronzi, Assya Trofimov, Brennan Nichyporuk, Justin Szeto, Nazanin Mohammadi Sepahvand, Edward Raff, Kanika Madan, Vikram Voleti, Samira Ebrahimi Kahou, Vincent Michalski, Tal Arbel, Chris Pal, Gael Varoquaux, and Pascal Vincent. Accounting for variance in machine learning benchmarks. In A. Smola, A. Dimakis, and I. Stoica (eds.), *Proceedings of Machine Learning and Systems*, volume 3, pp. 747–769, 2021. URL [https://proceedings.mlsys.org/paper\\_files/paper/2021/file/0184b0cd3cfb185989f858a1d9f5c1eb-Paper.pdf](https://proceedings.mlsys.org/paper_files/paper/2021/file/0184b0cd3cfb185989f858a1d9f5c1eb-Paper.pdf).
- Alexander Brenner, Felix Knispel, Florian P. Fischer, Peter Rossmanith, Yvonne Weber, Henner Koch, Rainer Röhrig, Julian Varghese, and Ekaterina Kutafina. Concept-based ai interpretability in physiological time-series data: Example of abnormality detection in electroencephalography. *Computer Methods and Programs in Biomedicine*, 257:108448, 2024. ISSN 0169-2607. doi: <https://doi.org/10.1016/j.cmpb.2024.108448>. URL <https://www.sciencedirect.com/science/article/pii/S0169260724004413>.
- Jeffrey W Britton, Lauren C Frey, Jennifer L Hopp, Pearce Korb, Mohamad Z Koubeissi, William E Lievens, Elia M Pestana-Knight, and Erk K St Louis. Electroencephalography (eeg): An introductory text and atlas of normal and abnormal findings in adults, children, and infants, 2016. URL <http://europepmc.org/books/NBK390354>.
- Jan Brogger, Tom Eichele, Eivind Aanestad, Henning Olberg, Ina Hjelland, and Harald Aurlien. Visual eeg reviewing times with score eeg. *Clinical neurophysiology practice*, 3:59–64, 2018.
- Zikun Cai, Tian jian Luo, and Xuan Cao. Multi-branch spatial-temporal-spectral convolutional neural networks for multi-task motor imagery eeg classification. *Biomedical Signal Processing and Control*, 93:106156, 2024. ISSN 1746-8094. doi: <https://doi.org/10.1016/j.bspc.2024.106156>. URL <https://www.sciencedirect.com/science/article/pii/S1746809424002143>.
- Zheng Yang Chin, Kai Keng Ang, Chuanchu Wang, Cuntai Guan, and Haihong Zhang. Multi-class filter bank common spatial pattern for four-class motor imagery bci. In *2009 Annual International Conference of the IEEE Engineering in Medicine and Biology Society*, pp. 571–574, 2009. doi: 10.1109/IEMBS.2009.5332383.
- P. Chrabąszcz. Neural architecture search. Master’s thesis, Albert Ludwig University Freiburg, 2018.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus), 2016. URL <https://arxiv.org/abs/1511.07289>.
- Mike X Cohen. *Analyzing neural time series data: theory and practice*. MIT press, 2014.
- Scott Cole and Bradley Voytek. Cycle-by-cycle analysis of neural oscillations. *Journal of neurophysiology*, 122(2):849–861, 2019.
- Alexander Craik, Yongtian He, and Jose L Contreras-Vidal. Deep learning for electroencephalogram (eeg) classification tasks: a review. *Journal of Neural Engineering*, 16(3):031001, apr 2019. doi: 10.1088/1741-2552/ab0ab5. URL <https://dx.doi.org/10.1088/1741-2552/ab0ab5>.
- Mohammad-Javad Darvishi-Bayazi, Mohammad Sajjad Ghaemi, Timothee Lesort, Md Rifat Arefin, Jocelyn Faubert, and Irina Rish. Amplifying pathological detection in eeg signaling pathways through cross-dataset transfer learning, 2023.
- Prabhu D Emmady and Arayamparambil C Anilkumar. *EEG abnormal waveforms*. StatPearls Publishing, Treasure Island (FL), 2023. URL <http://europepmc.org/books/NBK557655>.

- Taweesak Emsawas, Takashi Morita, Tsukasa Kimura, Ken-ichi Fukui, and Masayuki Numao. Multi-kernel temporal and spatial convolution for eeg-based emotion classification. *Sensors*, 22(21), 2022. ISSN 1424-8220. doi: 10.3390/s22218250. URL <https://www.mdpi.com/1424-8220/22/21/8250>.
- Denis A Engemann, Federico Raimondo, Jean-Rémi King, Benjamin Rohaut, Gilles Louppe, Frédéric Faugeras, Jitka Annen, Helena Cassol, Olivia Gosseries, Diego Fernandez-Slezak, et al. Robust EEG-based cross-site and cross-protocol classification of states of consciousness. *Brain*, 141(11):3179–3192, 10 2018. ISSN 0006-8950. doi: 10.1093/brain/awy251. URL <https://doi.org/10.1093/brain/awy251>.
- Stefan Falkner, Aaron Klein, and Frank Hutter. BOHB: Robust and efficient hyperparameter optimization at scale. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1437–1446. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/falkner18a.html>.
- Oliver Faust, Yuki Hagiwara, Tan Jen Hong, Oh Shu Lih, and U Rajendra Acharya. Deep learning for healthcare applications based on physiological signals: A review. *Computer Methods and Programs in Biomedicine*, 161:1–13, 2018. ISSN 0169-2607. doi: <https://doi.org/10.1016/j.cmpb.2018.04.005>. URL <https://www.sciencedirect.com/science/article/pii/S0169260718301226>.
- Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 25(5):845–869, 2014. doi: 10.1109/TNNLS.2013.2292894.
- Lukas A.W. Gemein, Robin T. Schirrmeister, Patryk Chrabąszcz, Daniel Wilson, Joschka Boedecker, Andreas Schulze-Bonhage, Frank Hutter, and Tonio Ball. Machine-learning-based diagnostics of eeg pathology. *NeuroImage*, 220:117021, 2020. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2020.117021>. URL <https://www.sciencedirect.com/science/article/pii/S1053811920305073>.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterton (eds.), *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pp. 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR. URL <https://proceedings.mlr.press/v9/glorot10a.html>.
- A. Gramfort, D. Strohmeier, J. Haueisen, M.S. Hämläinen, and M. Kowalski. Time-frequency mixed-norm estimates: Sparse m/eeg imaging with non-stationary source activations. *NeuroImage*, 70:410–422, 2013. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2012.12.051>. URL <https://www.sciencedirect.com/science/article/pii/S1053811912012372>.
- Shahram Hanifi, Andrea Cammarono, and Hossein Zare-Behtash. Advanced hyperparameter optimization of deep learning models for wind power prediction. *Renewable Energy*, 221:119700, 2024. ISSN 0960-1481. doi: <https://doi.org/10.1016/j.renene.2023.119700>. URL <https://www.sciencedirect.com/science/article/pii/S0960148123016154>.
- Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023. URL <https://arxiv.org/abs/1606.08415>.
- Matthias Hoppe. Eeg-befundung einschließlich darstellung des normalen eeg. *Das Neurophysiologie-Labor*, 40(1):14–43, 2018. ISSN 1439-4847. doi: <https://doi.org/10.1016/j.neulab.2017.11.002>. URL <https://www.sciencedirect.com/science/article/pii/S1439484717300571>. EEG Befunden.
- Edward E. Houfek and Robert J. Ellingson. On the reliability of clinical EEG interpretation. *The Journal of Nervous and Mental Disease*, 128(5):425–437, 1959.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. doi: <https://doi.org/10.48550/arXiv.1704.04861>.

- Thorir Mar Ingolfsson, Michael Hersche, Xiaying Wang, Nobuaki Kobayashi, Lukas Cavigelli, and Luca Benini. Eeg-tcnet: An accurate temporal convolutional network for embedded motor-imagery brain-machine interfaces. In *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pp. 2958–2965, 2020. doi: 10.1109/SMC42975.2020.9283028.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift, 2015. URL <https://arxiv.org/abs/1502.03167>.
- Md Kafiul Islam, Amir Rastegarnia, and Zhi Yang. Methods for artifact detection and removal from scalp eeg: A review. *Neurophysiologie Clinique/Clinical Neurophysiology*, 46(4-5):287–305, 2016.
- Mainak Jas, Denis A. Engemann, Yousra Bekhti, Federico Raimondo, and Alexandre Gramfort. Autoreject: Automated artifact rejection for meg and eeg data. *NeuroImage*, 159:417–429, 2017. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2017.06.030>. URL <https://www.sciencedirect.com/science/article/pii/S1053811917305013>.
- Ziyu Jia, Youfang Lin, Jing Wang, Kaixin Yang, Tianhang Liu, and Xinwang Zhang. Mmcnn: A multi-branch multi-scale convolutional neural network for motor imagery classification. In Frank Hutter, Kristian Kersting, Jefrey Lijffijt, and Isabel Valera (eds.), *Machine Learning and Knowledge Discovery in Databases*, pp. 736–751, Cham, 2021. Springer International Publishing. ISBN 978-3-030-67664-3.
- Wei-Bang Jiang, Li-Ming Zhao, and Bao-Liang Lu. Large brain model for learning generic representations with tremendous eeg data in bci, 2024. URL <https://arxiv.org/abs/2405.18765>.
- Yoonji Joo, Eun Namgung, Hyeonseok Jeong, Ilhyang Kang, Jinsol Kim, Sohyun Oh, In Kyoony Lyoo, Sujung Yoon, and Jaeuk Hwang. Brain age prediction using combined deep convolutional neural network and multi-layer perceptron algorithms. *Scientific Reports*, 13(1):22388, 2023.
- Nick Kane, Jayant Acharya, Sandor Beniczky, Luis Caboclo, Simon Finnigan, Peter W. Kaplan, Hiroshi Shibasaki, Ronit Pressler, and Michel J.A.M. van Putten. A revised glossary of terms most commonly used by clinical electroencephalographers and updated proposal for the report format of the eeg findings. revision 2017. *Clinical Neurophysiology Practice*, 2:170–185, 2017. ISSN 2467-981X. doi: <https://doi.org/10.1016/j.cnp.2017.07.002>. URL <https://www.sciencedirect.com/science/article/pii/S2467981X17300215>.
- Angus Kenny, Tapabrata Ray, and Hemant Singh. A framework for design optimization across multiple concepts. *Scientific Reports*, 14(1):7858, 2024.
- Hassan Aqeel Khan, Rahat Ul Ain, Awais Mehmood Kamboh, Hammad Tanveer Butt, Saima Shafait, Wasim Alamgir, Didier Stricker, and Faisal Shafait. The nmt scalp eeg dataset: An open-source annotated dataset of healthy and pathological eeg recordings for predictive modeling. *Frontiers in Neuroscience*, 15, 2022. doi: 10.3389/fnins.2021.755817. URL <https://www.frontiersin.org/article/10.3389/fnins.2021.755817>.
- Ann-Kathrin Kiessner, Robin T. Schirrmeister, Lukas A.W. Gemein, Joschka Boedecker, and Tonio Ball. An extended clinical eeg dataset with 15,300 automatically labelled recordings for pathology decoding. *NeuroImage: Clinical*, 39:103482, 2023. ISSN 2213-1582. doi: <https://doi.org/10.1016/j.nicl.2023.103482>. URL <https://www.sciencedirect.com/science/article/pii/S2213158223001730>.
- Ann-Kathrin Kiessner, Robin T. Schirrmeister, Joschka Boedecker, and Tonio Ball. Reaching the ceiling? empirical scaling behaviour for deep eeg pathology classification. *Computers in Biology and Medicine*, 178: 108681, 2024. ISSN 0010-4825. doi: <https://doi.org/10.1016/j.combiomed.2024.108681>. URL <https://www.sciencedirect.com/science/article/pii/S0010482524007662>.
- Min-jae Kim, Chul Youn Young, and Paik Joonki. Deep learning-based eeg analysis to classify normal, mild cognitive impairment, and dementia: Algorithms and dataset. *NeuroImage*, 272:120054, 2023. ISSN 1053-8119. doi: <https://doi.org/10.1016/j.neuroimage.2023.120054>. URL <https://www.sciencedirect.com/science/article/pii/S1053811923002008>.

- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- Elnaz Lashgari, Dehua Liang, and Uri Maoz. Data augmentation for deep-learning-based electroencephalography. *Journal of Neuroscience Methods*, 346:108885, 2020. ISSN 0165-0270. doi: <https://doi.org/10.1016/j.jneumeth.2020.108885>. URL <https://www.sciencedirect.com/science/article/pii/S0165027020303083>.
- Imene Latreche, Sihem Slatnia, Okba Kazar, and Saad Harous. An optimized deep hybrid learning for multi-channel eeg-based driver drowsiness detection. *Biomedical Signal Processing and Control*, 99:106881, 2025. ISSN 1746-8094. doi: <https://doi.org/10.1016/j.bspc.2024.106881>. URL <https://www.sciencedirect.com/science/article/pii/S174680942400939X>.
- Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and Brent J Lance. EEGNet: a compact convolutional neural network for EEG-based brain-computer interfaces. *Journal of neural engineering*, 15(5):056013, 2018. doi: [10.1088/1741-2552/aace8c](https://doi.org/10.1088/1741-2552/aace8c). URL <https://doi.org/10.1088/1741-2552/aace8c>.
- Tiago Lemos, Luiz Felipe Campos, Afrânio Melo, Nayher Clavijo, Rafael Soares, Maurício Câmara, Thiago Feital, Thiago Anzai, and José Carlos Pinto. Echo state network based soft sensor for monitoring and fault detection of industrial processes. *Computers & Chemical Engineering*, 155:107512, 2021. ISSN 0098-1354. doi: <https://doi.org/10.1016/j.compchemeng.2021.107512>. URL <https://www.sciencedirect.com/science/article/pii/S0098135421002908>.
- Chusak Limotai, Nittaya Phayaph, Boonyavee Pattanasilp, Jeerawan Mokklaew, and Natlada Limotai. Effects of antiepileptic drugs on electroencephalography (eeg): Insights and applicability. *Epilepsy & Behavior*, 110:107161, 2020. ISSN 1525-5050. doi: <https://doi.org/10.1016/j.yebeh.2020.107161>. URL <https://www.sciencedirect.com/science/article/pii/S1525505020303401>.
- Shiwei Liu, Wenwen Yue, Zhiqing Guo, and Liejun Wang. Multi-branch cnn and grouping cascade attention for medical image classification. *Scientific Reports*, 14(1):15013, 2024.
- Tianjun Liu and Deling Yang. A densely connected multi-branch 3d convolutional neural network for motor imagery eeg decoding. *Brain Sciences*, 11(2), 2021. ISSN 2076-3425. doi: [10.3390/brainsci11020197](https://doi.org/10.3390/brainsci11020197). URL <https://www.mdpi.com/2076-3425/11/2/197>.
- Xiaoguang Liu, Shicheng Xiong, Xiaodong Wang, Tie Liang, Hongrui Wang, and Xiuling Liu. A compact multi-branch 1d convolutional neural network for eeg-based motor imagery classification. *Biomedical Signal Processing and Control*, 81:104456, 2023. ISSN 1746-8094. doi: <https://doi.org/10.1016/j.bspc.2022.104456>. URL <https://www.sciencedirect.com/science/article/pii/S1746809422009107>.
- Sebas Lopez, G Suarez, D Jungreis, I Obeid, and Joseph Picone. Automated identification of abnormal adult eegs. In *2015 IEEE signal processing in medicine and biology symposium (SPMB)*, pp. 1–5. IEEE, 2015.
- S. López de Diego. Automated interpretation of abnormal adult electroencephalography. Master’s thesis, Temple University, 2017. URL <http://hdl.handle.net/20.500.12613/1767>.
- Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. doi: [10.48550/arXiv.1608.03983](https://doi.org/10.48550/arXiv.1608.03983).
- Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in Adam. *arXiv preprint arXiv:1711.05101*, 2017. doi: [10.48550/arXiv.1711.05101](https://doi.org/10.48550/arXiv.1711.05101).
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- H. B. Mann and D. R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *The Annals of Mathematical Statistics*, 18(1):50–60, 1947.

- L. McInnes, J. Healy, and J. Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *ArXiv e-prints*, February 2018.
- John William Carey Medithe and Usha Rani Nelakuditi. Study of normal and abnormal eeg. In *2016 3rd International Conference on Advanced Computing and Communication Systems (ICACCS)*, volume 01, pp. 1–4, 2016. doi: 10.1109/ICACCS.2016.7586341.
- Ghulam Muhammad, M. Shamim Hossain, and Neeraj Kumar. Eeg-based pathology detection for home health monitoring. *IEEE Journal on Selected Areas in Communications*, 39(2):603–610, 2021. doi: 10.1109/JSAC.2020.3020654.
- David O. Nahmias and Kimberly L. Kontson. Easy perturbation eeg algorithm for spectral importance (easypeasi): A simple method to identify important spectral features of eeg in deep learning models. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*, pp. 2398–2406, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403289. URL <https://doi.org/10.1145/3394486.3403289>.
- David O Nahmias, Kimberly L Kontson, David A Soltysik, and Eugene F Civillico. Consistency of quantitative electroencephalography features in a large clinical data set. *Journal of neural engineering*, 16(6): 066044, 2019.
- Chetan S Nayak and Arayamparambil C Anilkumar. Eeg normal waveforms. statpearls. *Treasure Island, FL: StatPearls Publishing*. <http://www.ncbi.nlm.nih.gov>, 2020.
- I. Obeid and J. Picone. The Temple University Hospital EEG data corpus. *Frontiers in Neuroscience*, 10, 2016. doi: 10.3389/fnins.2016.00196. URL <https://www.frontiersin.org/articles/10.3389/fnins.2016.00196>.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Olivier Petit, Nicolas Thome, Clement Rambour, Loic Themyr, Toby Collins, and Luc Soler. U-net transformer: Self and cross attention for medical image segmentation. In Chunfeng Lian, Xiaohuan Cao, Islem Rekik, Xuanang Xu, and Pingkun Yan (eds.), *Machine Learning in Medical Imaging*, pp. 267–276, Cham, 2021. Springer International Publishing. ISBN 978-3-030-87589-3.
- David Picard. Torch.manual\_seed (3407) is all you need: On the influence of random seeds in deep learning architectures for computer vision, 2023. URL <https://arxiv.org/abs/2109.08203>.
- Martyna Poziomska, Marian Dvoglalo, Przemysław Olbratowski, Paweł Niedbalski, Paweł Ogniewski, Joanna Zych, Jacek Rogala, and Jarosław Żygierewicz. Quantity versus diversity: Influence of data on detecting eeg pathology with advanced ml models, 2024. URL <https://arxiv.org/abs/2411.17709>.
- D. Merlin Praveena, D. Angelin Sarah, and S. Thomas George and. Deep learning techniques for eeg signal applications – a review. *IETE Journal of Research*, 68(4):3030–3037, 2022. doi: 10.1080/03772063.2020.1749143. URL <https://doi.org/10.1080/03772063.2020.1749143>.
- Anichur Rahman, Tanoy Debnath, Dipanjali Kundu, Md Saikat Islam Khan, Airin Afroj Aishi, Sadia Sazzad, Mohammad Sayduzzaman, and Shahab S Band. Machine learning and deep learning-based approach in smart healthcare: Recent advances, applications, challenges and opportunities. *AIMS Public Health*, 11(1):58, 2024.
- Mouad Riyad, Mohammed Khalil, and Abdellah Adib. Incep-eegnet: a convnet for motor imagery decoding. In *Image and Signal Processing: 9th International Conference, ICISP 2020, Marrakesh, Morocco, June 4–6, 2020, Proceedings 9*, pp. 103–111. Springer, 2020.
- Stephen W. Rose, J. Kiffin Penry, Billy G. White, and Susumu Sato. Reliability and Validity of Visual EEG Assessment in Third Grade Children. *Clinical Electroencephalography*, 4(4):197–205, 1973. doi: 10.1177/155005947300400405. URL <https://doi.org/10.1177/155005947300400405>.

- Subhrajit Roy, Isabell Kiral-Kornek, and Stefan Harrer. Deep learning enabled automatic abnormal eeg identification. In *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pp. 2756–2759, 2018. doi: 10.1109/EMBC.2018.8512756.
- Subhrajit Roy, Isabell Kiral-Kornek, and Stefan Harrer. Chrononet: A Deep Recurrent Neural Network for Abnormal EEG Identification. In David Riaño, Szymon Wilk, and Annette ten Teije (eds.), *Artificial Intelligence in Medicine*, pp. 47–56. Springer International Publishing, 2019a.
- Yannick Roy, Hubert Banville, Isabela Albuquerque, Alexandre Gramfort, Tiago H Falk, and Jocelyn Faubert. Deep learning-based electroencephalography analysis: a systematic review. *Journal of Neural Engineering*, 16(5):051001, aug 2019b. doi: 10.1088/1741-2552/ab260c. URL <https://dx.doi.org/10.1088/1741-2552/ab260c>.
- Zeynel A. Samak, Philip Clatworthy, and Majid Mirmehdi. Transop: Transformer-based multimodal classification for stroke treatment outcome prediction. In *2023 IEEE 20th International Symposium on Biomedical Imaging (ISBI)*, pp. 1–5, 2023. doi: 10.1109/ISBI53787.2023.10230576.
- Michiel Schinkel, Frank C Bennis, Anneroo W Boerman, W Joost Wiersinga, and Prabath WB Nanayakkara. Embracing cohort heterogeneity in clinical machine learning development: a step toward generalizable models. *Scientific reports*, 13(1):8363, 2023.
- R. T. Schirrmeister, L. A. W. Gemein, K. Eggersperger, F. Hutter, and T. Ball. Deep learning with convolutional neural networks for decoding and visualization of EEG pathology. *arXiv preprint arXiv:1708.08012*, 2017a.
- Robin Tibor Schirrmeister, Jost Tobias Springenberg, Lukas Dominique Josef Fiederer, Martin Glasstetter, Katharina Eggersperger, Michael Tangermann, Frank Hutter, Wolfram Burgard, and Tonio Ball. Deep learning with convolutional neural networks for eeg decoding and visualization. *Human Brain Mapping*, aug 2017b. ISSN 1097-0193. doi: 10.1002/hbm.23730. URL <http://dx.doi.org/10.1002/hbm.23730>.
- Shashank Shekhar, Adesh Bansode, and Asif Salim. A comparative study of hyper-parameter optimization tools. In *2021 IEEE Asia-Pacific Conference on Computer Science and Data Engineering (CSDE)*, pp. 1–6, 2021. doi: 10.1109/CSDE53843.2021.9718485.
- Unmesh Shukla, Geetika Jain Saxena, Manish Kumar, Anil Singh Bafila, Amit Pundir, and Sanjeev Singh. An improved decision support system for identification of abnormal eeg signals using a 1d convolutional neural network and savitzky-golay filtering. *IEEE Access*, 9:163492–163503, 2021. doi: 10.1109/ACCESS.2021.3133326.
- Hafza Ayesha Siddiqa, Zhenning Tang, Yan Xu, Laishuan Wang, Muhammad Irfan, Saadullah Farooq Abbasi, Anum Nawaz, Chen Chen, and Wei Chen. Single-channel eeg data analysis using a multi-branch cnn for neonatal sleep staging. *IEEE Access*, 12:29910–29925, 2024. doi: 10.1109/ACCESS.2024.3365570.
- S J M Smith. Eeg in the diagnosis, classification, and management of patients with epilepsy. *Journal of Neurology, Neurosurgery & Psychiatry*, 76(suppl 2):ii2–ii7, 2005. ISSN 0022-3050. doi: 10.1136/jnnp.2005.069245. URL [https://jnnp.bmj.com/content/76/suppl\\_2/ii2](https://jnnp.bmj.com/content/76/suppl_2/ii2).
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1): 1929–1958, 2014.
- Jens B Stephansen, Alexander N Olesen, Mads Olsen, Aditya Ambati, Eileen B Leary, Hyatt E Moore, Oscar Carrillo, Ling Lin, Fang Han, Han Yan, et al. Neural network analysis of sleep stages enables efficient diagnosis of narcolepsy. *Nature communications*, 9(1):5229, 2018.
- Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pp. 843–852, 2017.

- Haoqi Sun, Eyal Kimchi, Oluwaseun Akeju, Sunil B Nagaraj, Lauren M McClain, David W Zhou, Emily Boyle, Wei-Long Zheng, Wendong Ge, and M Brandon Westover. Automated tracking of level of consciousness and delirium in critical illness using deep learning. *NPJ Digital Medicine*, 2(1):1–8, 2019.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1–9, 2015. doi: 10.1109/CVPR.2015.7298594.
- Alon Talmor and Jonathan Berant. Multiqa: An empirical investigation of generalization and transfer in reading comprehension. *arXiv preprint arXiv:1905.13453*, 2019. URL <https://doi.org/10.48550/arXiv.1905.13453>.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6105–6114. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/tan19a.html>.
- IV Tatum and O William. *Handbook of EEG interpretation*. Springer Publishing Company, 2021.
- Juan Terven, Diana M. Cordova-Esparza, Alfonso Ramirez-Pedraza, Edgar A. Chavez-Urbiola, and Julio A. Romero-Gonzalez. Loss functions and metrics in deep learning, 2024. URL <https://arxiv.org/abs/2307.02694>.
- Rahul Thapa, Bryan He, Magnus Ruud Kjaer, Hyatt Moore IV, Gauri Ganjoo, Emmanuel Mignot, and James Y. Zou. SleepFM: Multi-modal representation learning for sleep across ECG, EEG and respiratory signals. In *AAAI 2024 Spring Symposium on Clinical Foundation Models*, 2024. URL <https://openreview.net/forum?id=cDXtscWCKC>.
- KG Van Leeuwen, H Sun, M Tabaeizadeh, AF Struck, MJAM Van Putten, and MB Westover. Detecting abnormal electroencephalograms using deep convolutional networks. *Clinical Neurophysiology*, 130(1): 77–84, 2019. doi: <https://10.1016/j.clinph.2018.10.012>. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6309707/>.
- D. Western, T. Weber, R. Kandasamy, F. May, S. Taylor, Y. Zhu, and L. Canham. Automatic report-based labelling of clinical eegs for classifier training. In *2021 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*, pp. 1–6, 2021. doi: 10.1109/SPMB52430.2021.9672295.
- Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training procedure in timm, 2021. URL <https://arxiv.org/abs/2110.00476>.
- Tao Wu, Xiangzeng Kong, Yiwen Wang, Xue Yang, Jingxuan Liu, and Jun Qi. Automatic classification of eeg signals via deep learning. In *2021 IEEE 19th International Conference on Industrial Informatics (INDIN)*, pp. 1–6, 2021. doi: 10.1109/INDIN45523.2021.9557473.
- Yuxin Wu and Kaiming He. Group normalization, 2018. URL <https://arxiv.org/abs/1803.08494>.
- Fei Yan, Zekai Guo, Abdullah M Ilyasu, and Kaoru Hirota. Multi-branch convolutional neural network with cross-attention mechanism for emotion recognition. *Scientific Reports*, 15(1):3976, 2025.
- Chaoqi Yang, M Westover, and Jimeng Sun. Biot: Biosignal transformer for cross-data learning in the wild. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 78240–78260. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/f6b30f3e2dd9cb53bbf2024402d02295-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/f6b30f3e2dd9cb53bbf2024402d02295-Paper-Conference.pdf).
- Lie Yang, Yonghao Song, Xueyu Jia, Ke Ma, and Longhan Xie. Two-branch 3d convolutional neural network for motor imagery eeg decoding. *Journal of Neural Engineering*, 18(4):0460c7, aug 2021. doi: 10.1088/1741-2552/ac17d6. URL <https://dx.doi.org/10.1088/1741-2552/ac17d6>.

- Özal Yildirim, Ulas Baran Baloglu, and U Rajendra Acharya. A deep convolutional neural network model for automated identification of abnormal eeg signals. *Neural Computing and Applications*, 32:15857–15868, 2020. doi: <https://doi.org/10.1007/s00521-018-3889-z>.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. URL <https://doi.org/10.48550/arXiv.1605.07146>.
- Aston Zhang, Zachary C. Lipton, Mu Li, and Alexander J. Smola. *Dive into Deep Learning*. Cambridge University Press, 2023a. <https://D2L.ai>.
- Hao Zhang, Qing-Qi Zhou, He Chen, Xiao-Qing Hu, Wei-Guang Li, Yang Bai, Jun-Xia Han, Yao Wang, Zhen-Hu Liang, Dan Chen, et al. The applied principles of eeg analysis methods in neuroscience and clinical neurology. *Military Medical Research*, 10(1):67, 2023b.
- Hangyu Zhu, Laishuan Wang, Ning Shen, Yonglin Wu, Shu Feng, Yan Xu, Chen Chen, and Wei 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*, 31: 2195–2204, 2023. doi: 10.1109/TNSRE.2023.3266876.

## A Appendix

### A.1 Related Work

In this section, we briefly summarise the related work on deep learning methods for general EEG pathology classification that have been successfully applied to the TUH Abnormal EEG Corpus (TUAB) (López de Diego, 2017) or TUH Abnormal Expansion Balanced EEG Corpus (TUABEXB) (Kiessner et al., 2023) (see Table 2 and Table 3 in Section 3.5). In recent years, a substantial amount of research has explored various approaches to general EEG pathology classification, reporting accuracies that lie within a narrow range of 81% to 87%. However, the majority of previous work is based on the TUAB (see Table 2 in Section 3.5). In this regard, Schirrmester et al. (2017a) introduced and applied two CNNs: Deep4Net and ShallowNet, which were originally proposed for motor imagery and then adapted for EEG pathology classification, achieving accuracies of 85.40% and 84.50%, respectively. These models<sup>4</sup> have been successfully reused or reimplemented in other studies (Darvishi-Bayazi et al., 2023; Gemein et al., 2020; Khan et al., 2022; Kiessner et al., 2023; 2024; Van Leeuwen et al., 2019; Western et al., 2021; Wu et al., 2021). Since then, several different architectures have been used to classify general EEG pathology. For instance, Gemein et al. (2020) employed two additional networks: a compact CNN called EEGNet and a temporal convolutional neural network (TCN). The EEGNet has been proposed by (Lawhern et al., 2018) for EEG-based BCIs using depthwise and separable convolutions. The TCN was introduced by (Bai et al., 2018) for the sequence modelling task and has then been optimised for EEG classification with a neural architecture search (Chrabaszcz, 2018; Gemein et al., 2020). Other studies have employed different CNNs, such as 1-D CNNs (Shukla et al., 2021; Yildirim et al., 2020), or hybrid models that combine CNN and RNN (Roy et al., 2018; 2019a), an Inception-Residual CNN (Wu et al., 2021), or an XceptionTime model (Brenner et al., 2024). Recently, a few approaches have explored the potential of a transformer-based foundation model for EEG pathology classification. The Large Brain Model (LaBraM) model (Jiang et al., 2024) (369M parameters) and the Biosignal Transformer (BIOT) model (Yang et al., 2023) (3.2M parameters) were pre-trained on a combination of different datasets using a cross-dataset training and evaluated on the TUAB dataset. Moreover, Roy et al. (2019a) proposed a deep 1D convolutional gated recurrent neural network, called ChronoNet, with exponentially varying filter sizes in the Conv1D layers. While ChronoNet achieved one of the highest accuracies (86.57%), its open-source re-implementation achieved a lower performance (81%) (Khan et al., 2022). Similarly, despite reusing the original implementation of some of these CNNs, studies reported slightly different results from the original report. The reasons for these differences include variations in the preprocessing of EEG data, such as different input lengths, different versions of the TUAB or different datasets, as well as differences in hardware

<sup>4</sup>are available in Braindecode, a deep learning toolbox for EEG, which is available for download at <https://github.com/TNLFreiburg/braindecode>.



and training strategies. On the TUAB, Wu et al. (2021) reported the highest accuracy of 87.10% using a model called IRCNN. Overall, the reported approaches achieved a mean specificity of above 84% and a mean sensitivity of about 75-80% on the TUAB. Most models show a lower sensitivity than specificity, indicating that they are more likely to classify pathological examples as non-pathological than vice versa. Baseline results on the TUABEXB dataset have been reported using Deep4Net, ShallowNet, TCN and EEGNet (see Table 3 in Section 3.5). These four CNNs achieved typical EEG pathology classification accuracies on TUABEXB, with slightly lower accuracies, ranging approximately from 82% to 86%. Consistent with previous studies on the TUAB, the models mostly achieved a lower sensitivity (72-82%) than specificity (83-96%) on the TUABEXB.

Although multi-branch, multi-scale, or parallel architectures (Altuwaijri et al., 2022; Belwafi et al., 2017; Ingolfsson et al., 2020; Jia et al., 2021; Riyad et al., 2020; Szegedy et al., 2015; Zhang et al., 2023a) are now beginning to gain traction in various EEG classification tasks, including motor imagery (Altuwaijri et al., 2022; Cai et al., 2024; Jia et al., 2021; Liu & Yang, 2021; Liu et al., 2023; Yang et al., 2021), emotion recognition (Emsawas et al., 2022; Yan et al., 2025), and neonatal sleep staging (Siddiqi et al., 2024; Zhu et al., 2023), only a few attempts have been made to assess the effectiveness of using CNNs with a set of three convolution scales for general EEG pathology classification. For example, Roy et al. (2019a) has considered the importance of convolutional layers with multiple filters of exponentially varying sizes (2, 4, and 8) to extract and combine features from different time scales. Inspired by the core idea of inception (Szegedy et al., 2015; Zhang et al., 2023a), Wu et al. (2021) employed stacked convolution kernels of different, yet small, scales (1, 3, and 5) in parallel. Recently, Brenner et al. (2024) employed an XceptionTime model, in which each module applies a set of three depthwise separable convolutions with different kernel sizes (11, 21 and 41). These studies have reported improvements in performance, with accuracies ranging from 85.10% to 87.10%. However, these results were based on the TUAB, and no attempts have been made to evaluate the performance of multi-scale CNNs on a larger, more heterogeneous dataset. In addition, smaller kernel sizes were preferred due to lower computational costs (Emsawas et al., 2022), which, however, tend to extract shorter temporal patterns from faster frequency bands (Cohen, 2014; Jia et al., 2021). In contrast, larger kernel sizes require higher computational costs, while they can learn long-term temporal patterns from slow frequency bands (Cohen, 2014; Jia et al., 2021). To date, research has not yet determined whether using raw EEG signals as input for a multi-branch CNN with multi-scale convolutions, which combine both long-term and short-term temporal patterns by employing kernel sizes based on clinically relevant frequency bands, is effective for classifying general EEG pathology on larger, heterogeneous datasets.

## A.2 Additional Details about the Datasets

In this section, we provide additional details on the TUH Abnormal EEG Corpus (TUAB)<sup>5</sup> (López de Diego, 2017) and the TUH Abnormal Expansion Balanced EEG Corpus (TUABEXB)<sup>6</sup> (Kiessner et al., 2023) that we used in our experiments. The two datasets are publicly available subsets of the Temple University Hospital EEG Corpus (TUEG)<sup>7</sup> (Obeid & Picone, 2016) for general EEG pathology classification, consisting of EEG recordings labelled as non-pathological or pathological. The TUEG comprises 69,582 clinical EEG recordings from 26,873 EEG sessions involving 15,001 patients over a 15-year period. Details on the number of TUAB recordings and patients are shown in Table 6, while Table 7 gives details of the number of recordings and patients in the TUABEXB dataset. For training, we used the concatenation of the TUAB training set and the TUABEXB training set, which we refer to as the TUH Abnormal Combined EEG Corpus (TUABCOMB). We evaluated our models on the predefined test set of both datasets, TUAB and TUABEXB, respectively (see Section 3.3). Details on the number of recordings and patients of the TUABCOMB are shown in Table 8. There is no overlap of patients between the training set of TUABCOMB and the evaluation sets of TUAB and TUABEXB.

<sup>5</sup>v2.0.0; available for download at [https://www.isip.piconepress.com/projects/tuh\\_eeg/html/downloads.shtml](https://www.isip.piconepress.com/projects/tuh_eeg/html/downloads.shtml).

<sup>6</sup>The EEG recordings are part of the TUH EEG Corpus and are publicly available after registration on the TUH EEG Corpus website ([isip.piconepress.com/projects/tuh\\_eeg/html/downloads.shtml](https://www.isip.piconepress.com/projects/tuh_eeg/html/downloads.shtml)), the corresponding pathology labels are available for download at [github.com/AKiessner/TUHAbnormal-Expansion-dataset](https://github.com/AKiessner/TUHAbnormal-Expansion-dataset).

<sup>7</sup>v1.1.0 and v1.2.0; available for download at [https://isip.piconepress.com/projects/tuh\\_eeg/downloads/tuh\\_eeg/](https://isip.piconepress.com/projects/tuh_eeg/downloads/tuh_eeg/).

Table 6: Number of recordings and patients in the TUAB dataset. For 54 patients in the training set, both non-pathological and pathological recordings are available. However, there is no overlap between patients in the training and evaluation sets. For more details on the dataset, see López de Diego (2017) and Obeid & Picone (2016).

<b>TUH Abnormal EEG Corpus (TUAB)</b>	<b>Training set</b>		<b>Evaluation set</b>	
	<b>Recordings</b>	<b>Patients</b>	<b>Recordings</b>	<b>Patients</b>
<b>Non-pathological</b>	1371	1237	150	148
<b>Pathological</b>	1,346	893	126	105
<b>Total</b>	2,717	2,130	276	253

Table 7: Number of recordings and patients in the TUABEXB dataset. There are 183 patients with both pathological and non-pathological recordings in the training set. The training and evaluation sets do not share recordings from the same patient. Details of the dataset and the labelling procedure can be found in Kiessner et al. (2023).

<b>TUH Abnormal Expansion Balanced EEG Corpus (TUABEXB)</b>	<b>Training set</b>		<b>Evaluation set</b>	
	<b>Recordings</b>	<b>Patients</b>	<b>Recordings</b>	<b>Patients</b>
<b>Non-pathological</b>	4,015	3,253	447	392
<b>Pathological</b>	3,975	3,166	442	378
<b>Total</b>	7,990	6,419	889	770

### A.3 Design Choices and Hyperparameter Optimisation

For our proposed Multi-BK-Net architecture described in Section 2.1, we evaluated several design choices, including architecture hyperparameters, such as total number of temporal filters, filter length of the later convolution layer, strides and types of non-linearities, as well as algorithm hyperparameters, such as learning rate, weight decay and number of training epochs. In addition, we evaluated potential performance improvements by using intermediate normalisation by batch normalisation (Ioffe & Szegedy, 2015) or group normalisation (Wu & He, 2018) as well as the use of exponential linear units (ELU,  $f(x) = x$  for  $x > 0$  and  $f(x) = e^x - 1$  for  $x \leq 1$  (Clevert et al., 2016) or GELU ( $GELU(x) = xP(X \leq x) = x\Phi(x)$ ) (Hendrycks & Gimpel, 2023) as activation function. The hyperparameters used and their possible values, which construct the search space, are listed in Table 9.

To identify the optimal set of the CNN hyperparameters (Table 9), we employed the automatic hyperparameter optimisation framework Optuna (Akiba et al., 2019a), as it has been successfully used to efficiently tune the hyperparameters of CNNs (Al-Ja’afreh et al., 2023; Hanifi et al., 2024; Latreche et al., 2025; Shekhar et al., 2021). Optuna implements Sequential Model-Based Optimisation (SMBO) to efficiently search for the optimal hyperparameters by building a probabilistic model of the objective function (Akiba et al., 2019a; Bergstra et al., 2011; Lemos et al., 2021). To search for the hyperparameter space and maximise the objective function, we used the Tree-structured Parzen Estimator (TPE) algorithm (Bergstra et al., 2011; 2013). We used multivariate TPE rather than independent TPE because multivariate TPE is reported to outperform independent TPE by finding better solutions faster and by better handling problems where there is an interaction between variables. The optimisation approach is described in detail in Bergstra et al. (2011; 2013); Falkner et al. (2018) and Kenny et al. (2024). In medical diagnosis, identifying all potentially pathological EEG recordings is more important than reducing false-positive results. Therefore, two objective measures — validation accuracy and validation sensitivity — were considered in this study. The multi-objective function for hyperparameter optimisation is defined as maximising the mean validation accuracy and mean validation sensitivity of the 5-fold cross-validation method. For each set of hyperparameters, we performed 5-fold cross-validation on the TUABCOMB training data using StratifiedGroupKFold, available in the Scikit-learn library (Pedregosa et al., 2011), with the constraint that patients were non-overlapping between splits (80% for training and 20% for validation). The splits were made with respect to the patients to avoid using record-

Table 8: Number of recordings and patients in the TUABCOMB training set dataset. There are 237 patients with both pathological and non-pathological recordings. The training set of the TUABCOMB and the evaluation sets of the TUAB and TUABEXB do not share the same patients.

TUABCOMB	Training set	
	Recordings	Patients
Non-pathological	5,386	4,490
Pathological	5,321	4,009
Total	10,707	8,549

Table 9: Configuration spaces considered in the search for the Multi-BK-Net architecture. We sampled until 100 configuration trials were completed to find the best setting. Learning rate and weight decay were defined as `trial.suggest_float(log=True)`.

Parameter	Config. Space
Total number of temporal conv filters	[20,25,30,35,40,45,50,55,60,65,70,75,80]
Normalisation	[batch, group]
Activation functions	[ELU, GELU]
Pooling mode first block	[mean, max]
Pooling mode other blocks	[mean, max]
Forth conv-pooling-block	[True, False]
Forth conv-pooling-block broader	[True, False]
Dropout	[0.4 - 0.6]
Filter length conv blocks	[10,15,20]
Input window size	6000
Weighted loss factor pathological	[1,2,3,4]
Optimiser	[AdamW]
Optimiser beta1	[0.5,0.9]
Learning rate	[ $1e^{-5}$ - $1e^{-1}$ ]
Weight decay	[ $1e^{-5}$ - $1e^{-1}$ ]
Batch size	[16,32,64]
Number of epochs	[30 - 105]
Number of channels	21

ings from the same patient in different folds. We additionally shuffled the data before splitting to ensure that there were no chronologically ordered splits, similar to the predefined training and evaluation sets in the TUAB and TUABEXB datasets. Experiments were conducted with a time budget of 45 hours per fold for each configuration run on a single fold. Runs that exceeded the time limit were pruned after completing the first fold. Runs that crashed (e.g., network configurations that did not fit in GPU memory) were also pruned. In addition, runs with a validation accuracy or validation sensitivity of less than 75% at the first fold were pruned to speed up optimisation. We ran the optimisation process until a total number of 100 trials were completed, where we jointly optimised parameters and architecture layouts. For the top 10 trials, we repeated the 5-fold cross-validation ten times and computed the mean validation accuracy and mean validation sensitivity, averaging across folds and then across runs. Finally, the best model configuration, with the highest mean accuracy and mean sensitivity, was used to train the model on the full TUABCOMB training set and then evaluated on the unseen, predefined evaluation sets from the TUAB and TUABEXB (see Section 3.3 for more details). The final hyperparameters are listed in Table 1 in Section 2.1.

#### A.4 Additional Details on the Baseline Architectures

In this section, we provide additional details on the baseline architectures, the list of hyperparameters for training and details on the training procedure. We additionally included five architectures implemented

Table 10: Optimised hyperparameters of the four single-scale CNNs that we used as a baseline (Gemein et al., 2020; Schirrmeister et al., 2017a; Kiessner et al., 2023; Chrabąszcz, 2018). No: Number of. For Deep4Net, ShallowNet and EEGNet, the parameter final\_conv\_length was set to 'auto' to enable trial-wise training. In case of TCN, an AdaptiveAvgPool1d layer was added to the final block.

Hyperparameter	Deep4Net	ShallowNet	TCN	EEGNet
No. input channels	21	21	21	21
Input time length	6000	6000	6000	6000
No. start filters	25	40	n.a.	n.a.
No. filters	n.a.	n.a.	55	n.a.
F1 (temporal filter)	n.a.	n.a.	n.a.	8
D depth multiplier	n.a.	n.a.	n.a.	2
F2 (pointwise filter)	n.a.	n.a.	n.a.	16
Final conv. length	67	394	n.a.	187
Channel factor	2	n.a.	n.a.	n.a.
Stride before pool	True	n.a.	n.a.	n.a.
Initial learning rate	0.01	0.000625	0.0011261049710243193	0.001
Weight decay	0.0005	0	5.83730537673086e-07	0
Dropout	0.5	0.5	0.05270154233150525	0.25
L2 decay	n.a.	n.a.	1.7491630095065614e-08	n.a.
Gradient clip	n.a.	n.a.	0.25	n.a.
Batch size	64	64	64	64
No. epochs	35	35	35	35
<b>Total No. parameters</b>	303,452	66,242	456,502	7,426

in Braindecode as a baseline: Deep4Net (Schirrmeister et al., 2017b;a), ShallowNet (Schirrmeister et al., 2017b;a), TCN (Bai et al., 2018; Chrabąszcz, 2018; Gemein et al., 2020), EEGNet Gemein et al. (2020); Lawhern et al. (2018) and ChronoNet Roy et al. (2019a). These architectures have been utilised in several studies and have demonstrated high performance on the TUAB and TUABEXB datasets across various versions of the TUAB and different preprocessing steps. The architectures are shown in Table 11 and Table 12. For a more detailed explanation of each architecture and its optimised hyperparameters, please refer to the corresponding studies. To ensure fairness in evaluation, we apply the same data preprocessing steps and training strategy to all architectures. While we have employed a trial-wise training strategy in this work, previous studies have mainly used a cropped training strategy (Gemein et al., 2020; Kiessner et al., 2023; 2024; Schirrmeister et al., 2017a). To adapt the architectures to our framework, we used their implemented approach for trial-wise training (Schirrmeister et al., 2017b). All models were trained on the TUABCOMB training set for 35 epochs using the trial-wise training strategy (Schirrmeister et al., 2017b). For all models, a batch size of 64 was used. The model parameters of Deep4Net, ShallowNet, TCN, and EEGNet were optimised using the AdamW optimiser (Loshchilov & Hutter, 2017). For more details on the hyperparameters used for training, see Table 10. For optimisation of the ChronoNet, we use the Adam optimiser (Kingma & Ba, 2017) with a learning rate of 0.001 and the binary cross-entropy loss. The learning rates for the gradient and weight decay updates were scheduled using cosine annealing (Loshchilov & Hutter, 2016), and we refrained from learning rate restarts.

## A.5 Additional Details on the Evaluation Metrics

In this section, we provide additional details on the evaluation metrics and significance tests used in this study. In our experiments, we evaluated the classification performance of our proposed Multi-BK-Net and comparison models using five performance evaluation metrics: accuracy, balanced accuracy, sensitivity, specificity and F2-score. These measures were determined as shown in Table 13. As both evaluation sets are slightly imbalanced (see Table 6 and Table 7), we also reported the balanced accuracy (BACC). Since our task is to classify EEG pathology, and thus no EEG recording containing pathological EEG patterns should be classified as non-pathological, it is important to reduce the false negative rates. False positives, on the

other hand, can be excluded by clinicians through a more detailed diagnostic process. As minimising false negatives is the main concern in this task, we additionally report the F2-score, which prioritises the model’s ability to identify pathological EEG recordings (recall) over the model’s ability to identify positive instances while minimising false positives.

## A.6 Detailed Results

### A.6.1 Performance Comparison to Baseline Approaches

The classification performance of the Multi-BK-Net and all baseline models on the TUAB is shown in Figure 6. On the TUAB evaluation set, Multi-BK-Net achieved a mean balanced accuracy of 87.36%, indicating the best classification performance, which was also statistically significantly better than the baseline models ( $p < 0.001$ , Mann-Whitney U test). Compared to four of the baseline architectures, our Multi-BK-Net achieved a slightly lower specificity, which was not statistically significantly different, except for the difference with Deep4Net ( $p < 0.001$ , Mann-Whitney U test). Figure 7 compares the balanced accuracy and sensitivity of our proposed Multi-BK-Net and all baseline architectures on the TUABEXB evaluation set. Again, the Multi-BK-Net outperformed the baseline models. In particular, the Multi-BK-Net achieved a higher mean balanced accuracy of 86.99% on the dedicated test set, which was also statistically significantly different from the baseline architectures ( $p < 0.001$ , Mann-Whitney U test). Although three of the baseline models achieved a slightly higher specificity than Multi-BK-Net (89.73%), only the difference between Multi-BK-Net and Deep4Net was statistically significant ( $p < 0.01$ ). Conversely, Multi-BK-Net achieved a statistically significantly higher mean specificity than ChronoNet ( $p < 0.05$ , Mann-Whitney U test).

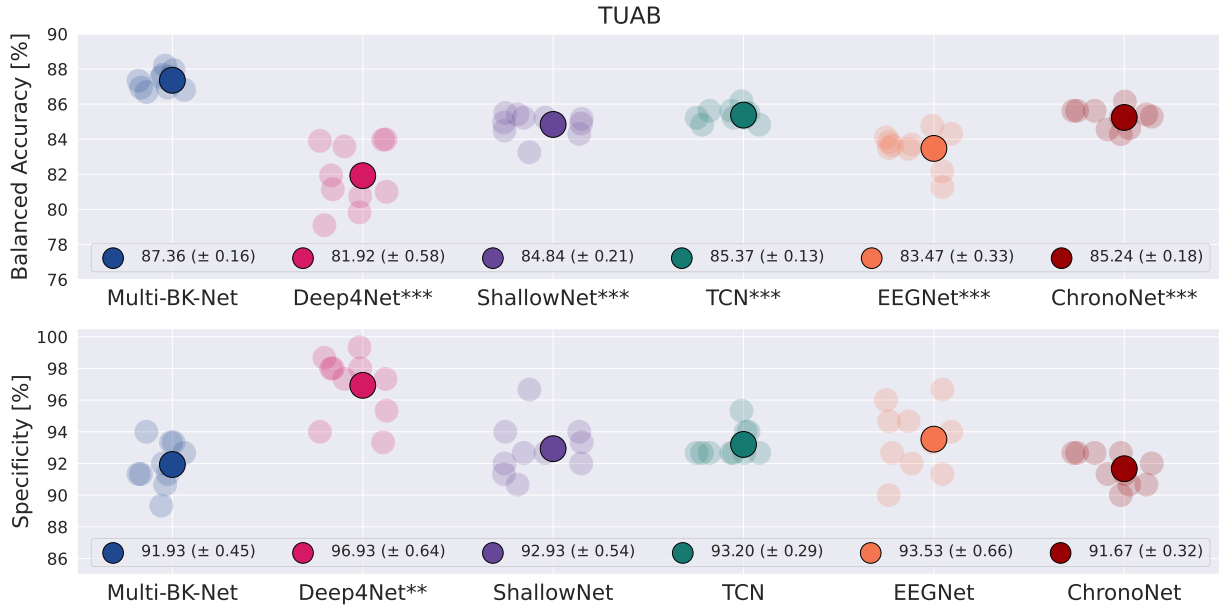


Figure 6: Performance comparison between our proposed Multi-BK-Net and five baseline architectures on the predefined TUAB evaluation set. Each transparent marker represents the performance of a single run, and each larger, bold symbol represents the mean performance score averaged across ten independent runs. The mean standard error is given in parentheses. Stars indicate statistically significant differences in performance score between the corresponding baseline architecture and the Multi-BK-Net (Bonferroni-corrected two-sided Mann-Whitney U test,  $p < 0.05$ : \*,  $p < 0.01$ : \*\*,  $p < 0.001$ : \*\*\*).

### A.6.2 Performance Comparison with Previous Reported State-of-the-art Deep Learning Approaches

Several studies have employed various end-to-end deep learning methods for the task of general EEG pathology classification (see Section A.1). In this section, we provide more details on the comparison of our proposed

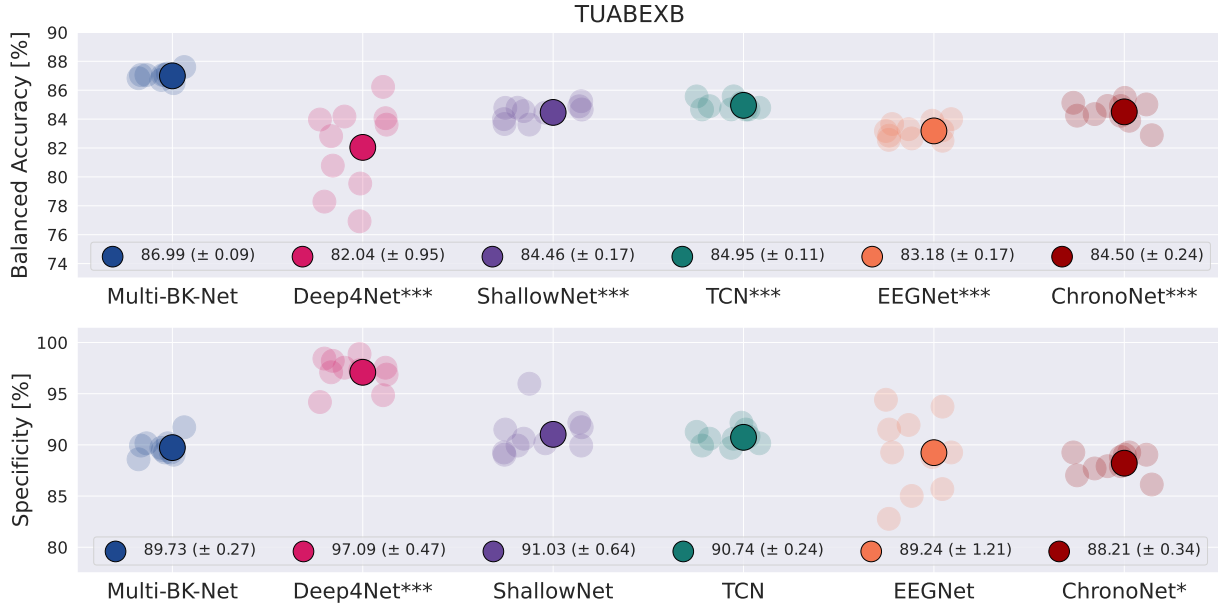


Figure 7: Performance comparison between our proposed Multi-BK-Net and all baseline architectures on the predefined TUABEXB evaluation set. Conventions as in Figure 6.

methods to previous work. For the selection of a previously published work on end-to-end deep learning models for general EEG pathology classification, we considered different architectures that have achieved competitive classification performance and match the following criteria: To ensure a fair comparison, the selected models adopt the raw or minimally preprocessed input from the same set of multiple standard scalp EEG electrodes and are evaluated on the predefined TUAB evaluation set or on the predefined TUABEXB evaluation set. Hence, publications that reported cross-validation results (Muhammad et al., 2021; Nahmias & Kontson, 2020; Poziomska et al., 2024), used single-channel data as input (Shukla et al., 2021; Yıldırım et al., 2020), or evaluated their methods on data from a non-publicly available dataset (Kim et al., 2023; Van Leeuwen et al., 2019), were excluded from our direct comparison. Table 2 and Table 3 in Section 3.5 summarise the comparison between our Multi-BK-Net and other methods on the TUAB and TUABEXB datasets, respectively. Perhaps the most clinically relevant finding is that Multi-BK-Net, with a mean sensitivity of 83.10%, outperforms all other approaches, reporting a mean sensitivity averaged over multiple runs, by at least 3.26% (see Table 2). Note that Brenner et al. (2024) reported a sensitivity of 85.70% for the XceptionTime model. However, the authors used only the first window of each recording as input during evaluation and reported single-run results, which is a major limitation of this approach. Due to the use of stochastic gradient-based algorithms with randomised weight initialisation, data ordering, and data augmentation during training, each independent training run produces a different network with better or worse performance than the average, with a variance of approximately 1% to 2% (Picard, 2023). Even if the variance between runs is not very large, there is a high probability that the performance of a single run will be an outlier, much better or much worse than the average (Picard, 2023). As a result, this variance can make it difficult to compare training configurations (Bouthillier et al., 2021; Picard, 2023). For example, comparing single runs, Multi-BK-Net achieved a sensitivity of 85.72% with accuracies of 88.41% and 87.68%, and a specificity of 90.67% and 89.33%, respectively, thus clearly outperforming the XceptionTime model. Nevertheless, to reduce stochasticity and improve comparability between training and model configurations, one should report the average of evaluation metrics over multiple runs (Wightman et al., 2021). Hence, a direct comparison of publications reporting single-run results would be unfair. Furthermore, comparing Table 2 with Table 3, it can also be seen that the mean accuracy of the Multi-BK-Net is slightly better on the TUAB than on the TUABEXB (+0.74%), while Multi-BK-Net achieved a higher mean sensitivity on TUABEXB than on TUAB (+1.15%). Interestingly, this result is consistent with previous observations (Kiessner et al., 2023), assuming that the differences are due to systematic differences in age, sex and patho-

logical feature distribution as well as different labelling methods (see Kiessner et al., 2023, for a more detailed discussion). Furthermore, we observed that Multi-BK-Net achieved a lower sensitivity than specificity on both evaluation sets, indicating a lower ratio of false positives than false negatives. This is also consistent with the results presented in previous literature (Gemein et al., 2020; Kiessner et al., 2023; 2024; López de Diego, 2017; Schirrmeister et al., 2017a; Western et al., 2021).

Overall, one of the main findings of our study is that the proposed model achieved mean accuracies of 87–88%. However, as noted and discussed in previous work (Engemann et al., 2018; Frénay & Verleysen, 2014; Gemein et al., 2020; Kiessner et al., 2024; Poziomska et al., 2024; Sun et al., 2019), there remains the possibility that this accuracy range may be the upper limit due to imperfect inter-rater agreement and the resulting label noise, with a model performance remaining below a perfect classification score (100%) and saturating around 88%. For the task of classifying EEG recordings as pathological or non-pathological, previous research has reported an inter-rater agreement of 86–88% between two neurologists (Houfek & Ellingson, 1959; Rose et al., 1973), while Beuchat et al. (2021) have found even lower mean inter-rater agreements of 82–86% between multiple EEG technologists and neurologists. Unfortunately, the differences in classification performance cannot be compared with the inter-rater agreement on these datasets used in this study for training and evaluation, as the labels of the TUAB and TUABEXB are mainly based on the impression of a single expert, as stated in the corresponding EEG report (for more details, see Kiessner et al., 2023; López de Diego, 2017; Obeid & Picone, 2016). Also, other studies on general EEG classification suffer from the same uncertainty. However, the need to include ratings from multiple experts has been recognised in other EEG classification tasks. For instance, Stephansen et al. (2018) compared the classification performance of their model with the inter-rater agreement of six raters for the task of classifying sleep stages. The authors found that the model’s accuracy increased from 76% (one rater) to 87% (consensus of six raters). Therefore, it can be assumed that the deep learning models for general EEG pathology also perform better than a single rater. It would thus be interesting to compare the classification performance of the proposed model with the inter-rater agreement of multiple raters. As the label noise remains the main limiting factor, it will be important to determine whether and to what extent label quality can be improved. For instance, using labels based on an ensemble of multiple raters would be a significant start to improving label quality. Additionally, using averaged ratings (so-called "fuzzy labels") for training could further enhance model performance, as they represent a probability of pathology, which is also helpful in identifying recordings that are ambiguous or more challenging to categorise, or contain more confounding patterns. Furthermore, as shown in Beuchat et al. (2021), multiple raters tend to achieve a higher mean sensitivity (87%) than mean specificity (75%), which is the opposite of machine learning approaches, which typically achieve higher specificity than sensitivity. Thus, including multiple raters as a baseline for training could increase the sensitivity of the classification model and also might further improve model performance, as was done for sleep stage classification. Moreover, a systematic optimisation of the labelling quality might even help the models overcome the current asymptotic limits of predictive accuracy (Darvishi-Bayazi et al., 2023; Kiessner et al., 2024; Poziomska et al., 2024), and hence lead to even better-performing models for the task of general EEG pathology classification.

### A.6.3 UMAP Visualisation of Learned Features

To further explain the proposed network, we visualised the learned features of the Multi-BK-Net before the classification layer using the UMAP (McInnes et al., 2018) method to map high-dimensional feature vectors into a two-dimensional UMAP space. This allowed the visual inspection of clustering and separation between pathological and non-pathological classes. For comparison, we also extracted and visualised the feature representations of the Deep4Net, as it is a single-scale CNN that is most similar in architecture to Multi-BK-Net, the TCN, which achieved the second-highest classification accuracy, and the ChronoNet, which uses a set of three convolutions. A UMAP comparison plot of the visualisation of the feature representation of the last layer before the classification layer is shown in Figure 8. In general, all models can form discernible clusters for pathological and non-pathological recordings, but with a noticeable amount of overlap. This is reflected in a reduced classification performance, which falls within a narrow range of 83% and 88%. As Multi-BK-Net shows the best classification performance, it also forms distinct and compact clusters for pathological and non-pathological recordings with minimal overlap. Smaller intra-class distances, i.e., tighter clustering, indicate a more consistent representation within each class, while larger inter-class distances, i.e.,

more distinct clusters, indicate that the model has learned features that effectively discriminate between the two classes. Compared to Multi-BK-Net, the three baseline models tend to form smaller and more dispersed clusters, possibly indicating greater within-class variability or a less precise representation of the pathological class. For example, Deep4Net forms less compact clusters, resulting in more dispersed clusters with larger intra-class distances. In addition, we can observe that pathological clusters are more dispersed along the component 1 axis, which may indicate greater within-class variability for pathological samples. This further suggests that these baseline models are less consistent in identifying pathological samples and are less capable of capturing the heterogeneity in pathological samples. The slightly lower effectiveness in capturing the characteristics of pathological samples may also reflect the lower sensitivity (high false negative rate) of the three baseline models compared to Multi-BK-Net. Additionally, when comparing the datasets (left vs. right column, Figure 4), we observe that all four models exhibit more distinct feature separation and better-defined clusters on the TUAB dataset than on the TUABEXB dataset. This suggests that the learned features are more discriminative for the TUAB dataset than the TUABEXB dataset. For all models, the distinction between pathological and non-pathological recordings was less pronounced on TUABEXB. Since better feature separation leads to improved classification performance, this may result in better model performance on TUAB than TUABEXB. A possible explanation for this could be that the TUABEXB is more heterogeneous compared to the TUAB (Kiessner et al., 2023), and as shown in previous work, this data heterogeneity can make EEG classification more difficult (Kiessner et al., 2023; Poziomska et al., 2024). Finally, the qualitative observations derived from the UMAP visualisations, which highlight the different degrees of feature separation achieved by different models, are consistent with quantitative performance metrics, such as accuracy and sensitivity. For instance, the distinct and compact clustering of data points, indicative of strong feature discrimination, observed for the Multi-BK-Net corresponds to higher accuracy and sensitivity scores compared to other models on both datasets. Conversely, the large intra-class distances seen in the visualisations for the pathological samples of the three baseline CNNs align with their lower sensitivity performance (67-81%). Similarly, models with poorer cluster separation in the UMAP visualisations, such as the Deep4Net, show comparable poorer quantitative performance (accuracy of 82-83%). This consistency between qualitative visualisations and quantitative results supports the validity of the UMAP analysis and provides a comprehensive understanding of model performance.

#### A.6.4 Ablation experiments

In this section, we provide additional details on the ablation experiments (Section 3.6) and present additional results. Compared to the five baseline architectures, our proposed Multi-BK-Net is significantly larger (1,038,683 parameters). For example, it is approximately twice the size of TCN (456,502 parameters), three times the size of Deep4Net (303,452 parameters), and eight times the size of ChronoNet (137,233 parameters). The Multi-BK-Net achieved an accuracy improvement of 2-5% over the baseline models, while sensitivity increased by more than 3.5-17.3%. Previous studies have shown that larger models perform better than smaller models (Darvishi-Bayazi et al., 2023; Howard et al., 2017; Kiessner et al., 2024; Talmor & Berant, 2019; Tan & Le, 2019; Sun et al., 2017; Zagoruyko & Komodakis, 2016). To verify that the superiority of the Multi-BK-Net over the baseline models was not due to the increased model size, we performed an ablation experiment. As the Deep4Net has the most similar architecture to our proposed Multi-BK-Net, we compared the performance of our approach with that of three larger versions of the Deep4Net (see Table 14). In particular, we increased the number of temporal filters in the Deep4Net to 35, which is the same number of filters used in the first convolution-pooling block of the Multi-BK-Net. We refer to this model as Deep4Net35, which comprises 579,182 parameters, approximately 55% of the size of Multi-BK-Net. We have also increased the number of temporal filters in the Deep4Net to 47 (Deep4Net47). This increases the number of trainable parameters to 1,026,482, which is approximately the relative size of 99% of the Multi-BK-Net. Finally, we created a variant called Deep4Net48, which contains 48 temporal filters and has 1,069,490 trainable parameters; thus, it is slightly larger than Multi-BK-Net.

To further validate the effectiveness of our multi-temporal kernel branches in the first convolution-pooling block of our proposed Multi-BK-Net, we also trained and evaluated different variants of our model, in which we removed four of the five branches while increasing the number of filters in the remaining branch from 7 to 35. Thus, the variant models consist of a single branch with 35 filters, with only one temporal kernel length in the first block. The kernel length of the variants is set to either 3, 7, 10, 13, 25 or 200. The Multi-BK-Net



uses kernel lengths of 3, 7, 13, 25 and 200, while the Deep4Net uses a kernel size of 10. We refer to these variants as STKSBNNet (Single-Temporal-Kernel Single-Branch Net). For example, STKSBNNet3 refers to the model variant that uses a temporal kernel length of 3. Furthermore, to investigate the effect of different branches, we included an additional variant with five different temporal kernels (3, 7, 13, 25, and 200), but only one branch was used instead of five. Thus, the features of the five temporal convolutional layers are concatenated directly before the spatial convolution. We refer to this model variant as MTKSBNNet (Multi-Temporal-Kernel Single-Branch Net). For more details on the size of the models used, see Table 14. We trained the STKSBNNet and MTKSBNNet variants as described in Section 2.1, and the Deep4Net variants as described in Section 3.2 and evaluated all variants as described in Section 3.3.

As shown in Table 4 in Section 3.6, the performance of Deep4Net increases as the number of filters increases. However, only Deep4Net47 obtained a statistically significantly higher mean F2-score and mean sensitivity on the TUABEXB set compared to Deep4Net ( $p < 0.05$ , one-sided Mann-Whitney U test). All other differences between the original Deep4Net and its larger variants were not statistically significant ( $p > 0.05$ , Mann-Whitney U test). Nevertheless, Multi-BK-Net outperforms all variants of Deep4Net statistically significantly on both datasets ( $p < 0.01$ , Mann-Whitney U test). This suggests that the performance improvement of Multi-BK-Net compared to Deep4Net is not due to increased model width or model size.

Moreover, Table 5 in Section 3.6 shows the performance of the variant models and the Multi-BK-Net on the TUAB and TUABEXB evaluation sets. It can be seen that when four branches and four of the kernel sizes are removed, there is a significant decrease in the results on both evaluation sets ( $p < 0.05$ , Mann-Whitney U test). By using a single kernel length, STKSBNets cannot extract accurate features, resulting in lower overall performance compared to Multi-BK-Net. This highlights the effectiveness of using multiple temporal kernel lengths. Furthermore, when comparing the performance of the STKSBNNet variants, it is evident that the optimal kernel size for STKSBNNet varies across different datasets. For TUAB, the mean accuracy is highest when the kernel size is 13, while the optimal kernel size for TUABEXB is 200. This observation is consistent with previous studies, which report that determining the optimal kernel length for EEG pathology classification is challenging due to subject and time differences (Jia et al., 2021). A possible solution is to use multiple kernels of different sizes, which also improves classification performance (Altuwaijri et al., 2022; Jia et al., 2021). Furthermore, the results show that the inclusion of multiple temporal kernels of different lengths, but with a single branch, improves performance compared to using a single kernel of a single length. This shows that multiple kernels can learn valuable features and improve classification performance. However, as shown in Table 5, incorporating multiple branches within the first block further improves performance on both datasets compared to using only a single branch (MTKSBCNN). When the four branches are removed, but multiple kernel lengths are used (MTKSBCNN), the mean F2-score and mean sensitivity decrease ( $p < 0.05$ , Mann-Whitney U test). This shows that the introduction of multiple input branches improves the identification of pathological samples. Overall, this ablation experiment demonstrates that incorporating both multiple temporal kernel lengths and multiple branches, as is done in Multi-BK-Net, is beneficial for classifying EEG pathology and significantly enhances the model’s performance.

Table 11: Architecture of the four single-scale, baseline CNNs adapted for trial-wise training and as implemented in Braindecode (Gemein et al., 2020; Schirrmeyer et al., 2017a; Kiessner et al., 2023; Chrabąszcz, 2018). E: Number of input electrodes. St: Stride.

Deep4Net	ShallowNet	TCN	EEGNet
25×Conv2D (10×1)	40×Conv2D (25×1)	55×Conv1D (16)	8×Conv2D (1×64)
25×Conv2D (1×E)	40×Conv2D (1×E)	WeightNorm	BatchNorm
BatchNorm	BatchNorm	Activation (ReLU)	16×Conv2D (E×1)
Activation (ELU)	Activation (Square)	Dropout2d (0.25)	BatchNorm
MaxPool (3×1)	MeanPool (75×1) St(15×1)	55×Conv1D (16)	Activation (ELU)
	Activation (Log)	WeightNorm	MeanPool (1×4) St(1×4)
	Dropout (0.5)	Activation (ReLU)	Dropout (0.25)
		Dropout2d (0.25)	
		55×Conv1D (1)	
		Activation (ReLU)	
Dropout (0.5)	2×Conv2D (394×1)	55×Conv1D (16)	16×Conv2D (1×16)
50×Conv2D (10×1) St(3×1)	LogSoftmax	WeightNorm	16×Conv2D (1×1)
BatchNorm		Activation (ReLU)	BatchNorm
Activation (ELU)		Dropout2d (0.25)	Activation (ELU)
MaxPool (3×1)		55×Conv1D (16)	MeanPool (1×8) St(1×8)
		WeightNorm	Dropout (0.25)
		Activation (ReLU)	
		Dropout2d (0.25)	
Dropout (0.5)		55×Conv1D (55×16)	2×Conv2D (1×187)
100×Conv2D (10×1) St(3×1)		WeightNorm	LogSoftmax
BatchNorm		Activation (ReLU)	
Activation (ELU)		Dropout2d (0.25)	
MaxPool (3×1)		55×Conv1D (16)	
		WeightNorm	
		Activation (ReLU)	
		Dropout2d (0.25)	
Dropout (0.5)		55×Conv1D (16)	
200×Conv2D (10×1) St(3×1)		WeightNorm	
BatchNorm		Activation (ReLU)	
Activation (ELU)		Dropout2d (0.25)	
MaxPool (3×1)		55×Conv1D (16)	
		WeightNorm	
		Activation (ReLU)	
		Dropout2d (0.25)	
2×Conv2D (67×1)		55×Conv1D (16)	
LogSoftmax		WeightNorm	
		Activation (ReLU)	
		Dropout2d (0.25)	
		55×Conv1D (16)	
		WeightNorm	
		Activation (ReLU)	
		Dropout2d (0.25)	
		Linear (55,2)	
		AdaptiveAvgPool1d	
		LogSoftmax	

Table 12: Architecture of the baseline architecture ChronoNet as proposed by Roy et al. (2019a). ChronoNet includes multiple filters of exponentially varying lengths in the 1D convolutional layers and dense connections within the GRU layers. Arrows represent skip connections. K: kernel\_size. st: stride. p: padding. h: hidden\_size

ChronoNet		
$32 \times \text{Conv1D}(k(2,), \text{st}(2,), p(1,))$	$32 \times \text{Conv1D}(k(4,), \text{st}(2,))$	$32 \times \text{Conv1D}(k(8,), \text{st}(2,), p(3,))$
	Filter Concat	
$32 \times \text{Conv1D}(k(2,), \text{st}(2,), p(1,))$	$32 \times \text{Conv1D}(k(4,), \text{st}(2,))$	$32 \times \text{Conv1D}(k(8,), \text{st}(2,), p(3,))$
	Filter Concat	
$32 \times \text{Conv1D}(k(2,), \text{st}(2,), p(1,))$	$32 \times \text{Conv1D}(k(4,), \text{st}(2,))$	$32 \times \text{Conv1D}(k(8,), \text{st}(2,), p(3,))$
	Filter Concat	
	GRU(h(32))	
	GRU(h(32))	
	Filter Concat	
	GRU(h(32))	
	Filter Concat	
	GRU(h(32))	
	Linear layer	
<b>Total No. parameters</b>	137,233	

Table 13: Performance evaluation metrics.  $TP$  is the number of examples that were correctly classified in the positive class,  $TN$  is the number of examples that were correctly classified in the negative class,  $FP$  is the number of examples that were incorrectly classified in the positive class,  $FN$  is the number of examples that were incorrectly classified in the negative class. F2-score:  $\beta = 2$  to emphasise the need to minimise false negatives, and  $Precision = \frac{TP}{TP+FP}$  and  $Recall = \frac{TP}{TP+FN}$ . Higher values of the F2-score indicate better performance in terms of recall.

Performance metric	Mathematical formula
Accuracy	$\frac{(TP+TN)}{(TP+FN+TN+FP)}$
Balanced accuracy	$\frac{1}{2}(\frac{TP}{TP+FN} + \frac{TN}{TN+FP})$
Sensitivity	$\frac{TP}{TP+FN}$
Specificity	$\frac{TN}{TN+FP}$
F2-score	$(1+\beta^2)(\frac{\frac{TP}{TP+FP} * \frac{TP}{TP+FN}}{(\beta^2 * \frac{TP}{TP+FP}) + \frac{TP}{TP+FN}})$

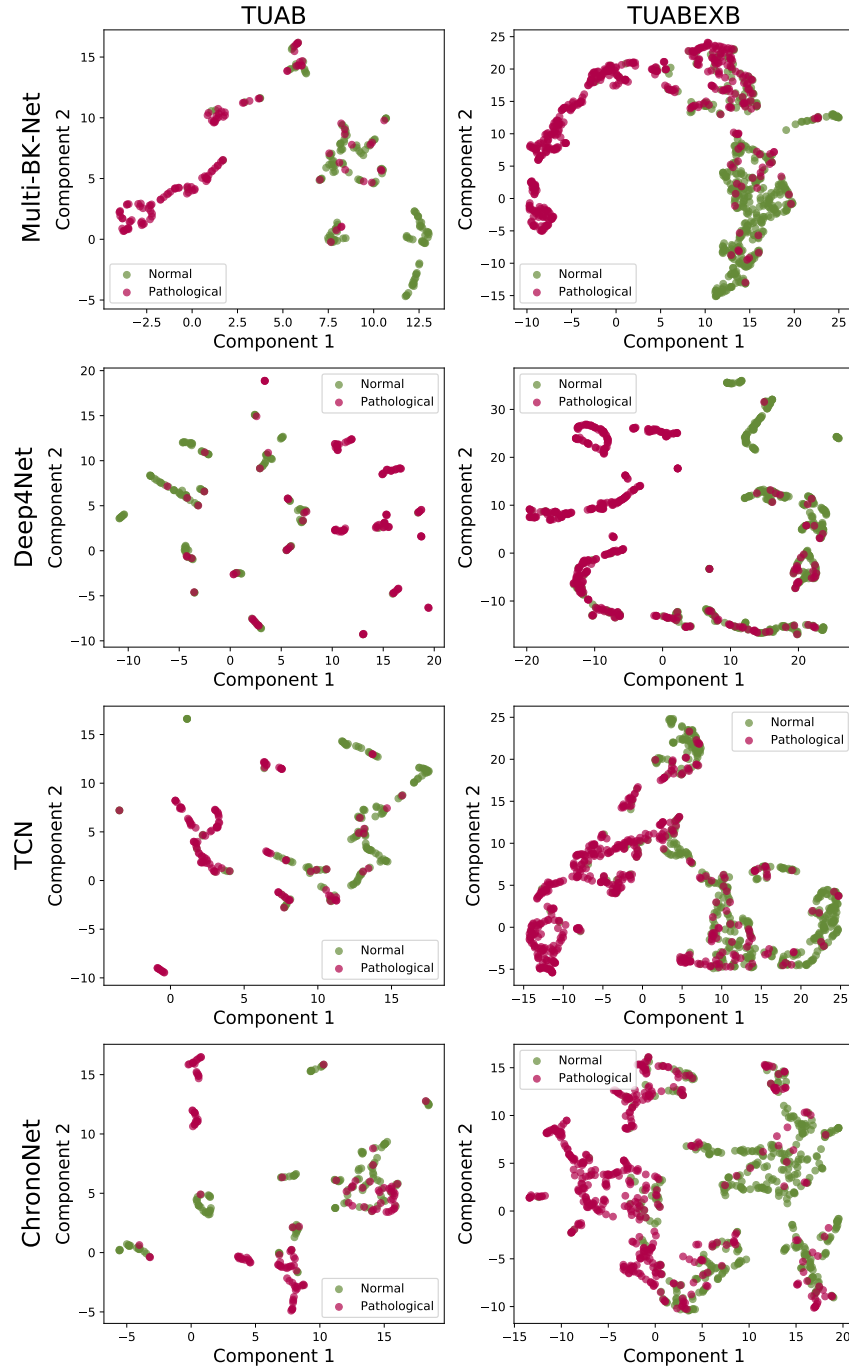


Figure 8: A comparison of UMAP visualisation of feature representations learned with Multi-BK-Net, Deep4Net, TCN and ChronoNet for pathological and non-pathological recordings on the predefined TUAB (left column) and the TUABEXB (right column) evaluation sets. Rows represent the UMAP visualisation of different models. Pink dots represent the pathological class, while green dots represent the non-pathological class.

Table 14: Size of the different architectures used in the ablation experiments. Model size is given as the number of trainable parameters.

Architecture	Model size
Deep4Net	303,452
Deep4Net35	579,182
Deep4Net47	1,026,482
Deep4Net48	1,069,490
STKSBN3	1,057,632
STKSBN7	1,057,772
STKSBN10	1,057,877
STKSBN13	1,057,982
STKSBN25	1,058,402
STKSBN200	1,064,527
<b>Multi-BK-Net</b>	<b>1,038,683</b>