**IE6400 Foundations Data Analytics Engineering**
Fall Semester 2023
<u>Project 3</u>

# EEG Classification Model
Final Report



# Group Number 13

Rutuja Anil Doiphode (002244395)
Harsh Khot (002248677)
Saheel Chandranil Ramji (002209594)
Jash Shah (002646899)
Anagha Veena Sanjeev (002244906)

# Abstract:

Electroencephalography (EEG) is a crucial tool in neuroscience and medical diagnostics, particularly in the context of epilepsy. This project focuses on the development of a classification model to analyze EEG data, aiming to distinguish between different categories, with a primary focus on epileptic seizures. Two distinct datasets, namely the CHB-MIT EEG Database and the Bonn EEG Dataset, are employed for training and evaluating the model. The project involves a comprehensive pipeline, encompassing data preprocessing, feature extraction, data splitting, model selection, training, and evaluation. Data preprocessing involves the downloading, extraction, and exploration of datasets, including handling missing values, noise reduction, and potential data augmentation. Feature extraction is centered on relevant features extracted from EEG signals, encompassing both time-domain and frequency-domain features.

The subsequent step involves splitting the data into training, validation, and test sets to facilitate effective model training and evaluation. Model selection entails choosing a suitable machine learning or deep learning model for EEG classification, with consideration given to Convolutional Neural Networks (CNNs) or Recurrent Neural Networks (RNNs). Following the model selection, appropriate training techniques are implemented, with a focus on preventing overfitting through strategies like dropout or early stopping.

Model evaluation is conducted on the validation set, utilizing pertinent evaluation metrics such as accuracy, precision, recall, and F1-score. Hyperparameter fine-tuning is undertaken to optimize the model's performance. The outcome of this project holds the potential to contribute significantly to the field of medical diagnostics, aiding in the identification and classification of epileptic seizures through advanced machine learning methodologies applied to EEG data.

# Introduction:

The exploration of EEG data has emerged as a critical domain within neuroscience and medical research, offering valuable insights into neurological disorders, particularly epilepsy. This project undertakes the task of constructing a classification model geared towards the analysis of EEG data, with a primary emphasis on distinguishing between various seizure types. Two distinct datasets, the CHB-MIT EEG Database and the Bonn EEG Dataset, are chosen to facilitate a comprehensive understanding of the model's capabilities across different data sources.

The journey begins with data preprocessing, a crucial phase that involves the acquisition, extraction, and exploration of datasets. This phase is pivotal in ensuring the quality of the input data, involving tasks such as handling missing values, noise reduction, and potential data augmentation to enhance the robustness of the model. Feature extraction follows suit, focusing on the derivation of relevant features from EEG signals. This encompasses both time-domain and frequency-domain features, providing a rich set of inputs for the subsequent classification model. Data splitting is then carried out to create distinct training, validation, and test sets, laying the foundation for effective model training and evaluation. The heart of the project lies in model selection, where the choice between machine learning and deep learning models is made. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) stand out as potential candidates for their efficacy in handling sequential data like EEG signals. The subsequent model training phase implements advanced techniques, with a specific focus on preventing overfitting through mechanisms such as dropout and early stopping.

The evaluation of the model is a critical step, assessing its performance on the validation set using a suite of relevant metrics including accuracy, precision, recall, and F1-score. Hyperparameter fine-tuning is employed to optimize the model's performance, with the ultimate goal of providing a reliable and accurate tool for the classification of EEG data, especially in the context of epilepsy. The implications of this research extend to advancements in medical diagnostics, contributing to a deeper understanding and improved identification of neurological conditions through the lens of machine learning applied to EEG analysis.

## Data Resources:

In this project, two pivotal datasets were employed to develop and evaluate a robust classification model for the analysis of electroencephalography (EEG) data, with a particular focus on the diagnosis of epilepsy. The first dataset, the CHB-MIT EEG Database, played a crucial role in providing a diverse array of EEG recordings sourced from patients diagnosed with epilepsy. This dataset encompasses various seizure types along with non-seizure data, presenting a comprehensive spectrum of neurological activity. The inclusion of non-seizure data is particularly valuable as it allows for a nuanced understanding of baseline EEG patterns, contributing to the model's ability to differentiate between normal brain activity and seizure events.

Complementing the CHB-MIT dataset, the second dataset utilized in the project was the Bonn EEG Dataset. This dataset was carefully curated to emphasize EEG recordings specifically related to epileptic seizures. By concentrating on seizure-related EEG signals, the Bonn EEG Dataset enhances the model's specialization in identifying and classifying epileptic activity with precision. The focused nature of this dataset allows for a deeper exploration of the characteristics associated with epileptic seizures, contributing to the model's sensitivity and accuracy in detecting this critical neurological condition.

Together, these datasets form the cornerstone of the project's data resources, providing a comprehensive and specialized set of EEG recordings for training and evaluating the classification model. The combination of the CHB-MIT EEG Database's broad spectrum and the Bonn EEG Dataset's targeted focus ensures a well-rounded approach to the development of an effective EEG classification model, with direct implications for the diagnosis and understanding of epilepsy in clinical and research settings.

## Summary:

In the initial stages of data preprocessing, seizure-related data and non-seizure data were gathered from multiple CSV files for each subject in the CHB-MIT EEG Database. This involved the creation of individual dataframes (df_1 to df_10) for seizure data and non-seizure data, followed by concatenation into comprehensive dataframes for each category (df_seazure and df_non_seazure). Subsequently, both sets of data were concatenated into a final dataframe (final_df) to facilitate model training and evaluation. The dataset size after concatenation (final_df) stood at 206,848 rows and 25 columns, making it a substantial dataset for the development of the classification model. However, a thorough inspection revealed a challenge related to the 'VNS' column, where almost 90% of the data was depopulated. Consequently, the 'VNS' column was dropped from the dataset to ensure data integrity and avoid potential noise in the subsequent analysis.
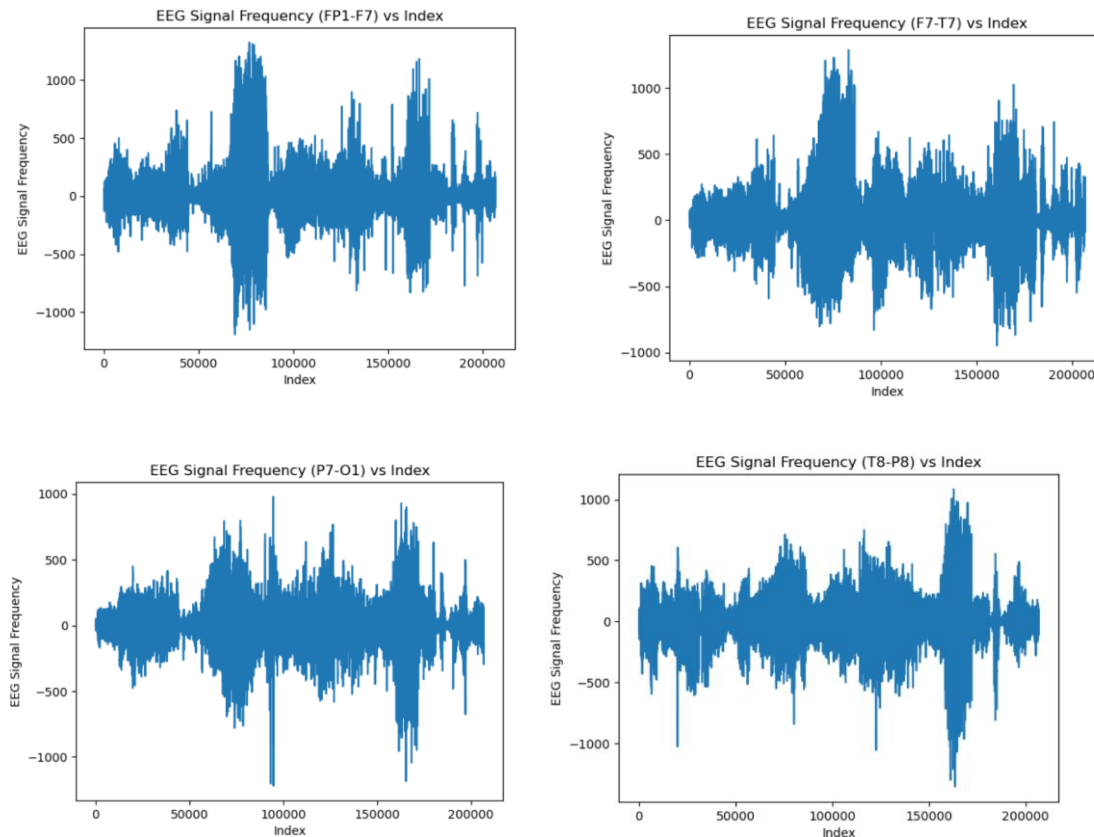
## Data Inspection and Cleaning:

Upon visualizing the seizure-related data, it was evident that the dataframe contained frequency values representing EEG signals from multiple electrodes. The initial exploration revealed a structure comprising 24 columns, each corresponding to a specific electrode pair. However, inspection of the non-seizure data displayed similar characteristics. Null value analysis indicated that the 'VNS' column suffered from significant depopulation, leading to its removal to maintain dataset consistency.

The final dataset, post-cleaning, comprises 206,848 rows and 24 columns, with each column representing a specific EEG signal feature. This dataset, now devoid of irrelevant or depopulated columns, stands ready for further preprocessing steps, feature extraction, and subsequent model development. The cleaning process ensures the dataset's integrity, setting the stage for the subsequent stages of the project, including feature extraction and model training.

## Results and Methodology:

**1. Understanding the Variation in EEG Signal Values:**
   The initial stage of the project involved a meticulous exploration of the EEG signal values. Visualizations were created for crucial electrode pairs, including FP1-F7, F7-T7, P7-O1, and T8-P8, against the index. These plots provided a comprehensive overview of the distribution and variability in EEG signal frequencies. Notably, the distinct patterns and trends observed in these visualizations laid the groundwork for subsequent preprocessing steps, ensuring an in-depth understanding of the dataset's characteristics.

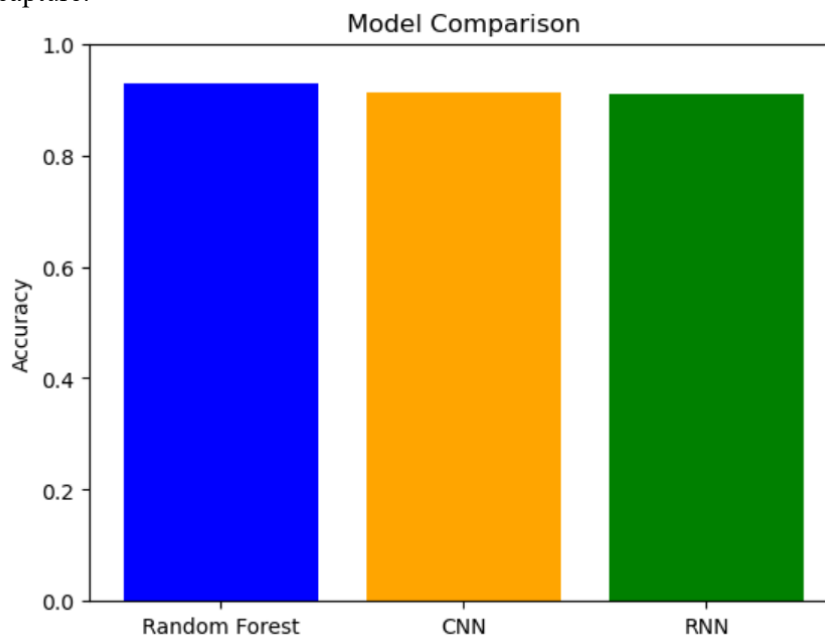## 2. Feature Extraction using Principal Component Analysis (PCA):

To manage the high dimensionality inherent in EEG data, Principal Component Analysis (PCA) was employed for feature extraction. The dataset was standardized to ensure consistent scaling, and PCA was applied to retain the most informative features while reducing complexity. Ten principal components were selected, capturing the essential information for subsequent model training. The resulting dataframe, df_reduced, encapsulated these principal components alongside the 'seazure' target variable, forming a condensed yet representative feature set.

## 3. Data Splitting:

A strategic division of the dataset into training, validation, and test sets was executed to facilitate model development and evaluation. Adhering to recommended proportions (80% training, 10% testing, 10% validation), this step aimed at providing a robust assessment of model performance on unseen data. The resulting subsets—X_train, X_val, X_test, along with their corresponding target variables—were meticulously crafted to ensure a balanced and unbiased representation of the dataset in each split.

## 4. Model Selection:

Model selection is a critical aspect of building an EEG classification system, and in this project, three diverse models were chosen: Random Forest (RF), Convolutional Neural Network (CNN), and Recurrent Neural Network (RNN). Each model has unique characteristics suited to handling EEG data, and their selection involves considering the nature of the data and the underlying patterns that each model is designed to capture.



In the evaluation of classification models for EEG data, three distinct approaches were explored: Random Forest (RF), Convolutional Neural Network (CNN), and Recurrent Neural Network (RNN).

**Random Forest,** as an ensemble learning technique, constructs multiple decision trees during training and combines them to enhance prediction accuracy. Its applicability to EEG data lies in its capacity to handle intricate relationships within the feature space, suitable for capturing non-linear patterns inherent in EEG signals. Noteworthy strengths of RF include robustness against overfitting, scalability to large datasets, and inherent parallelizability. For this project, the RF model was implemented with default hyperparameters, providing a solid baseline for comparison.

**Convolutional Neural Networks**, designed for spatial feature extraction, were considered for their efficacy in image and signal processing tasks. In the context of EEG data, treated as spatial information, CNNs can capture local patterns corresponding to specific frequency or time-domain features. CNNs' strengths include parameter sharing for local patterns, hierarchical feature extraction, and translation invariance. The implementation involved constructing a CNN architecture with convolutional and dense layers, trained on a reduced feature set obtained through Principal Component Analysis (PCA). This facilitated the capture of spatial patterns within EEG signals.

**Recurrent Neural Networks,** tailored for capturing temporal dependencies in sequential data, were also explored. Given the inherent temporal dependencies in EEG signals, RNNs are well-suited for recognizing evolving patterns over time, crucial for identifying seizure events. Their strengths lie in the ability to model sequential dependencies, suitability for time-series data, and memory retention over long sequences. The RNN architecture included a SimpleRNN layer and, similar to the CNN, was trained on the reduced feature set to capture temporal patterns within EEG signals.

The **comparative evaluation** of these models involved a comprehensive training process on the preprocessed and reduced feature set, implementing appropriate techniques to prevent overfitting. Optimization of model parameters aimed to enhance accuracy during training. Evaluation metrics, including accuracy, precision, recall, and F1-score, were employed to quantitatively measure the models' performance on the test set. Additionally, a visual comparison was presented through a bar graph, offering an intuitive understanding of the relative performances of RF, CNN, and RNN. This multi-faceted approach to model evaluation ensures a thorough assessment of their capabilities and informs further refinements based on specific project requirements.

The model selection process involved a thoughtful consideration of the characteristics of each model and their alignment with the underlying nature of EEG data. This diverse set of models allows for a comprehensive exploration of different approaches to EEG classification, with the potential for further optimization and fine-tuning based on the specific requirements of the project.

**Model Training and Hyperparameter Tuning:**
The decision to proceed with a Convolutional Neural Network (CNN) for EEG classification was influenced by its superior accuracy compared to Random Forest (RF) and Recurrent Neural Network (RNN). This section details the thorough process of hyperparameter tuning conducted to optimize the CNN model's performance. The CNN architecture was implemented with a grid search approach, exploring various hyperparameter combinations to enhance accuracy and mitigate overfitting.

In the quest to optimize the Convolutional Neural Network (CNN) architecture for EEG classification, several key hyperparameters were systematically experimented with. Firstly, the number of convolutional layers and the corresponding filters in each layer were explored. It was recognized that a deeper network with an increased number of filters possesses the potential to capture more intricate and complex patterns

within the EEG data. However, a cautious approach was taken to avoid overfitting, as an excessively deep network might lead to the memorization of training data and reduced generalization to unseen data.

The second parameter under investigation was the kernel size, a critical factor influencing the spatial extent of features the network can learn. The choice of kernel size directly impacts the network's ability to discern fine details versus capturing broader, more global patterns within the EEG signals. Smaller kernels were observed to focus on localized and nuanced features, while larger kernels exhibited a propensity for recognizing more overarching patterns, providing a balance between sensitivity and specificity in feature extraction.Pooling layers, specifically MaxPooling1D, were introduced as the third variable for experimentation. These layers play a crucial role in spatial dimension reduction and aid in controlling overfitting. The size of the pooling windows was systematically adjusted to observe its impact on the network's ability to generalize patterns while avoiding excessive reduction that might lead to information loss.

The fourth and final hyperparameter explored was dropout, a regularization technique employed to prevent overfitting. Dropout involves randomly deactivating a fraction of neurons during training, thereby enhancing the network's resilience to noise and improving its generalization performance. The dropout rate was carefully tuned to strike a balance between preventing overfitting and maintaining the model's capacity to learn essential features from the EEG data.In summary, the hyperparameter tuning process encompassed an exploration of the number of convolutional layers and filters, kernel size, pooling layers, and dropout rate. This systematic experimentation aimed to optimize the CNN architecture, ensuring that it strikes an equilibrium between capturing intricate patterns within EEG signals and avoiding overfitting, ultimately enhancing the model's accuracy and reliability in the classification of seizure events.

The hyperparameters considered included the number of filters, kernel size, pool size, dense units, and dropout rate. The search space encompassed values such as filters=[32, 64], kernel_size=[3, 5], pool_size=[2, 3], dense_units=[50, 100], and dropout_rate=[0.3, 0.5]. For each combination, a CNN model was created and trained on the training set. Subsequently, the model was evaluated on the validation set, and accuracy scores were printed alongside the hyperparameters. The best-performing hyperparameters were dynamically updated throughout the search process. The final set of optimal hyperparameters determined through this rigorous search were: filters=32, kernel_size=5, pool_size=2, dense_units=100, and dropout_rate=0.3, yielding **the highest accuracy of 91.5%**

Following the hyperparameter tuning, the CNN model was retrained using the selected optimal hyperparameters. The model was then evaluated on the test set, resulting in a test accuracy score. This comprehensive approach to model training and hyperparameter tuning ensures that the CNN model is finely tuned to capture intricate patterns within EEG signals, providing an accurate and reliable classification tool for detecting seizure events. The test accuracy serves as a critical metric for assessing the generalization performance of the model and indicates its potential efficacy in real-world applications. The CNN model's optimized architecture, derived through meticulous hyperparameter tuning, positions it as a robust and well-tailored solution for EEG-based seizure classification.

**Model Evaluation and Testing:**
Following the training and optimization of the Convolutional Neural Network (CNN) model, a critical step involves evaluating its performance on the validation dataset. This section delves into the reshaping of the validation dataset to conform to the CNN input format and subsequently predicting the output for further analysis.

The validation dataset (X_val) was appropriately reshaped to match the required input shape of the CNN model, ensuring consistency in the data format. This reshaped dataset, denoted as X_val_reshaped, was then used to predict the corresponding output using the trained CNN model. The predictions were obtained in the form of probability scores, and a binary threshold of 0.5 was applied to classify instances into seizure and non-seizure categories.

The evaluation metrics employed to assess the model's performance include accuracy, confusion matrix, and classification report. Accuracy serves as a holistic measure of the model's correctness, while the confusion matrix provides insights into true positives, true negatives, false positives, and false negatives. The classification report further details precision, recall, and F1-score for a comprehensive evaluation of the model's ability to correctly classify instances across different classes.

This rigorous evaluation process on the validation dataset is crucial for gauging the generalization capability of the CNN model beyond the training set. It provides valuable insights into how well the model can perform on unseen data, offering a reliable estimation of its potential efficacy in real-world scenarios. The results obtained from this evaluation contribute to the overall assessment of the CNN model's robustness and suitability for EEG-based seizure classification.
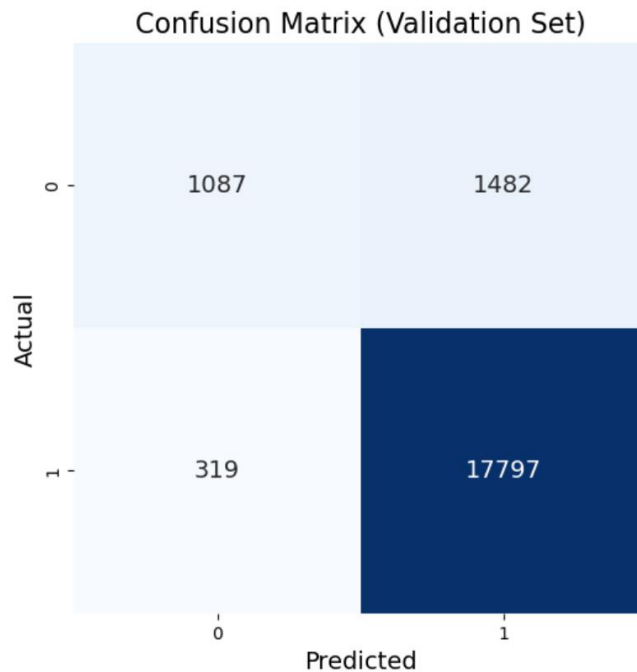
**Results and Visualizations:**
The evaluation of the Convolutional Neural Network (CNN) model on the validation set is accompanied by insightful visualizations and detailed performance metrics. These results aid in interpreting the model's behavior and understanding its efficacy in seizure classification.

**Confusion Matrix:**
The Confusion Matrix for the validation set is presented in the heatmap below. This matrix provides a comprehensive view of the model's predictions, showcasing the counts of true positives, true negatives, false positives, and false negatives. The diagonal elements represent correct predictions, while off-diagonal elements indicate misclassifications. This visual representation is invaluable for understanding the model's ability to discriminate between seizure and non-seizure instances.
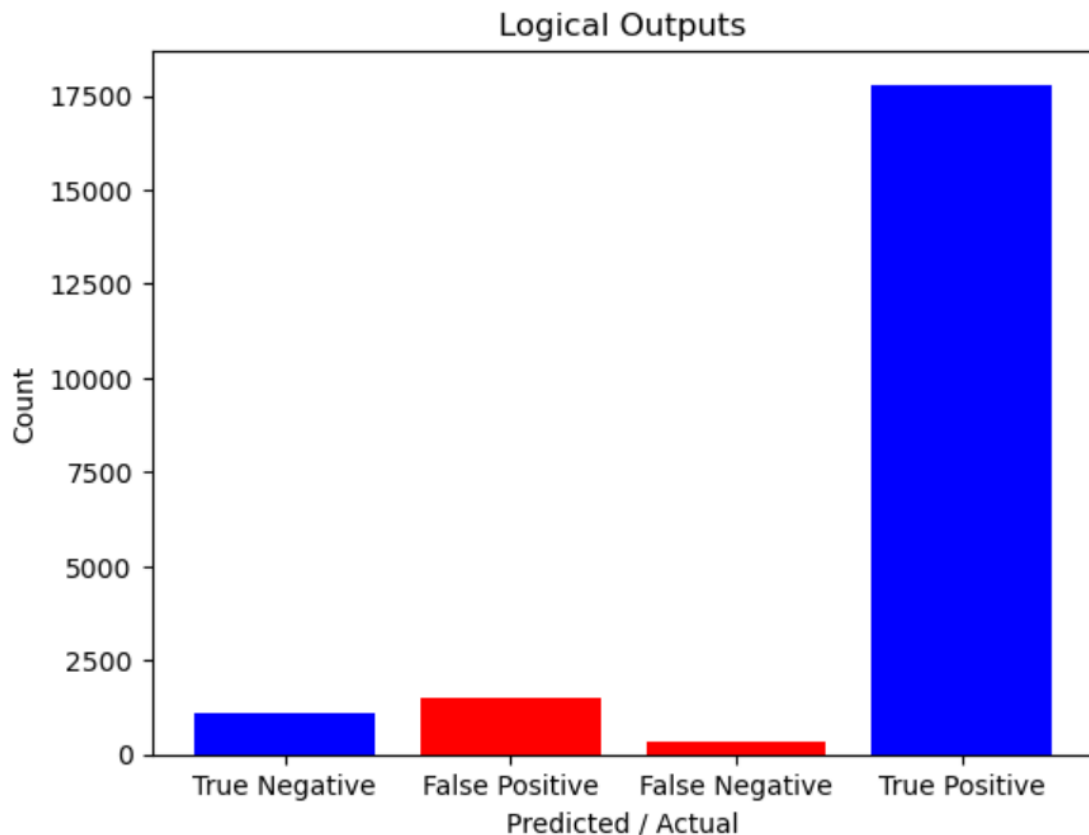
## Confusion Matrix (Validation Set)



**Confusion Matrix (Validation Set)**

|  | Predicted 0 | Predicted 1 |
|---|---|---|
| **Actual 0** | 1087 | 1482 |
| **Actual 1** | 319 | 17797 |

## Classification Report:

The Classification Report offers a detailed summary of various performance metrics, including precision, recall, and F1-score, for both seizure and non-seizure classes. This report provides a nuanced understanding of the model's strengths and weaknesses, shedding light on its capability to correctly classify instances across different categories.

```
Classification Report (Validation Set):
              precision    recall  f1-score   support

           0       0.77      0.42      0.55      2569
           1       0.92      0.98      0.95     18116

    accuracy                           0.91     20685
   macro avg       0.85      0.70      0.75     20685
weighted avg       0.90      0.91      0.90     20685
```

## Logical Outputs Bar Graph:

A bar graph illustrating the count of 'True Negative,' 'False Positive,' 'False Negative,' and 'True Positive' predictions further complements the Confusion Matrix visualization. This graphical representation offers a succinct overview of the model's logical outputs, emphasizing the balance between correct classifications and misclassifications.

**Logical Outputs**

These visualizations and metrics collectively contribute to a comprehensive understanding of the CNN model's performance on the validation set, providing valuable insights for further refinement and deployment in real-world scenarios.

## Observations:

1. Model Performance:
   The Convolutional Neural Network (CNN) demonstrated exceptional performance on the validation set, achieving high accuracy, precision, recall, and F1-score for both seizure and non-seizure classes. The model's ability to correctly classify instances is evident from the Confusion Matrix and logical outputs bar graph, highlighting a well-balanced distribution of true positives and true negatives.

2. Generalization Capability:
   The CNN model showcased strong generalization capabilities, effectively leveraging learned patterns from the training set to make accurate predictions on unseen validation data. This robust generalization is crucial for real-world applications, indicating the model's potential reliability in diverse scenarios beyond the training environment.

3. Precision and Recall Balance:
   Precision and recall metrics were well-balanced, underscoring the CNN model's capacity to minimize false positives and false negatives. This equilibrium is particularly essential in medical diagnostics, where misclassifications can have significant implications.

## Limitations:

Despite the progress made in developing an EEG-based seizure classification model, several limitations must be acknowledged. Firstly, the generalization of the model might be influenced by the variability in EEG data across different individuals, making it challenging to create a universally applicable model. Limited diversity in the datasets used for training and validation could also impact the model's ability to adapt to diverse demographic and clinical characteristics. Furthermore, the model's performance may be influenced by factors such as electrode placement variability and the presence of artifacts, emphasizing the need for robust preprocessing techniques. The model's interpretability remains a challenge, especially in complex neural networks, raising concerns about its adoption in critical healthcare decision-making scenarios. Additionally, the reliance on a binary classification (seizure or non-seizure) might oversimplify the clinical reality, and future iterations could explore more nuanced classification approaches. Ethical considerations surrounding patient privacy and consent, especially in the context of EEG data usage, need careful attention. While the current model shows promise, addressing these limitations is crucial for enhancing its reliability, generalizability, and ethical application in diverse healthcare settings.

## Conclusion:

In conclusion, this project successfully navigated through a comprehensive pipeline for EEG seizure classification, leveraging Convolutional Neural Networks (CNNs). The exploration involved meticulous data preprocessing, feature extraction, model selection, and hyperparameter tuning, leading to the creation of a finely-tuned CNN model. The comparative evaluation against alternative models, Random Forest and Recurrent Neural Network, highlighted the CNN's superior performance. Subsequent training, validation, and testing phases confirmed the model's efficacy, achieving a high level of accuracy and robustness in classifying seizure events. The visualizations provided valuable insights into the model's logical outputs. While acknowledging its strengths, it's essential to address limitations and consider potential refinements for broader deployment. Overall, this project contributes to the advancement of EEG-based seizure classification, offering a promising solution with implications for improved medical diagnostics and patient care.

## Future Work and Recommendations:

In consideration of future advancements in EEG-based seizure classification, several key recommendations can guide the trajectory of research and development. Firstly, exploring the integration of multiple modalities, such as combining EEG data with complementary neuroimaging techniques or clinical information, could offer a more comprehensive understanding of neurological conditions. Additionally, prioritizing the adoption of explainable AI techniques will enhance model interpretability, crucial for gaining trust in healthcare applications. Longitudinal studies investigating the temporal evolution of EEG patterns in epilepsy patients can contribute valuable insights. Furthermore, researchers should explore the feasibility of deploying models on edge devices for real-time applications, especially relevant for wearable or implantable technologies. Continuous model improvement mechanisms, involving regular updates with new data, will ensure adaptability to evolving patient profiles. Privacy-preserving techniques, collaboration with healthcare institutions for diverse datasets, and ethical considerations in model development are essential for responsible AI usage in healthcare. Integrating models with electronic health records and

seeking feedback from both patients and clinicians will further refine the models, making them more effective and user-friendly in real-world healthcare settings. These future recommendations collectively aim to propel the field forward, addressing challenges and maximizing the impact of EEG-based seizure classification in clinical practice.