

# A Root-Guided Random Forest Framework with Discriminative Voice Biomarker Selection for Parkinson's Disease Detection

Geetha Ramani R.<sup>1</sup>, Nandhitha K.<sup>2</sup>

Department of Information Technology, Anna University, Chennai, India.

E-mail: <sup>1</sup>rggeetha@yahoo.com, <sup>2</sup>nnandhitha50@gmail.com

## Abstract

Parkinson's Disease (PD) is a neurodegenerative disorder that severely impacts speech production by causing unstable phonation, articulation disorders, and prosody impairments. Early identification of such speech-specific PD symptoms enables timely diagnosis and treatment of the disease. The paper introduces a Root-Guided Random Forest Feature Selection (RGRFFS) framework for automatic detection of Parkinson's Disease using speech data. Samples of voices collected from the existing publicly available Parkinson's disease speech datasets have been preprocessed by performing noise reduction, voice activity detection, normalization, and signal segmentation. In total, 107 acoustic and spectral features were extracted, which include Mel-Frequency Cepstral Coefficients (MFCCs), measures of voice perturbations, formants, spectral features, and prosodic parameters, to capture speech features of Parkinsonian patients. To minimize redundant information and increase discriminative power of the set of features, Root-Guided Random Forest Feature Selection method was used for selection of vocal biomarkers ranked by Root Importance Score (RIS). As a result, 75 most informative features have been chosen and used for further classification and evaluation of classification accuracy. The achieved classification accuracy of 93.85% together with precision, recall, and F1-score values of 0.94 indicates that the selected feature subset allows reliable separation of samples of PD and control speech.

**Keywords:** Parkinson's Disease, Random Forest, Feature Selection, Voice Biomarkers, Root-Guided Learning.

## 1. Introduction

Parkinson's disease (PD) is a degenerative neurological condition that affects millions of people globally and is marked by motor dysfunction, cognitive impairment, and speech disorders [1], [2]. Among the many clinical manifestations of PD, speech disturbance has been found to be an early and common signature. The patients usually suffer from low voice loudness, monotonous voice pitch, difficulty in articulation, voice instability, and aberrant prosody, thereby making speech analysis a feasible method for early detection and monitoring of the disease [3], [4].

Modern advancements in machine learning, and signal processing techniques for speech signals can lead to an automated system that can detect the Parkinson disease through analyzing the acoustic features acquired from recorded speech samples [5], [6]. Many researchers have used feature extraction combined with classifiers to distinguish PD-specific abnormalities in speech. Since speech data have a massive number of acoustic features, a significant portion of those features are usually found to be redundant, have high correlations or little informativeness, leading to higher computational costs and a deterioration in classification accuracy [7], [8].

Most existing feature importance techniques compute global importance across the entire classifier's decision process and might not effectively reflect feature's importance at the decision formative process stages [9]. Therefore, developing a feature selection method that emphasizes features with better class separation and increased interpretability will be highly beneficial. In this research, the proposed Root-Guided Random Forest Feature Selection (RGRFFS) framework for the task of detecting Parkinson disease through speech. The feature's root-level discriminant is used to assign a Root Importance Score (RIS), which will select informative vocal biomarker candidates for PD detection. A comprehensive set of acoustic and spectral acoustic descriptors is extracted from speech samples, and more discriminative features were then selected for classification.

## 2. Literature Survey

Experimental evaluation showed the efficiency of the method to reduce redundant features, while maintaining a high discrimination rate. The result proves the usefulness of a root-guided feature selection for building an accurate and interpretable automated system to assist in detecting PD. Speech analysis is considered a non-invasive, promising alternative approach in detecting Parkinson disease (PD) in recent years by investigating and analyzing various clinical features of the disease from speech samples [3].

In the research done by many researchers [3] [10] [2] [4] [9], many machine learning algorithms such as Support Vector Machines, Naive Bayes, and Decision Trees, etc have been proposed and evaluated for this classification task by extracting the following biomarkers: Mel-Frequency Cepstral Coefficients (MFCCs), jitter, shimmer, pitch, prosodic acoustic features and spectral features. These models reported outstanding accuracies in PD classification and interpretability which can be clinically relevant. Deep representations has attracted much attention. Recent studies have explored self-supervised representation learning for speech analysis by extracting high-level embeddings from unlabeled audio data. These approaches leverage pre-trained speech models to capture discriminative vocal characteristics and have demonstrated promising performance in voice disorder detection tasks [10]. The comparison highlights the growing adoption of representation learning techniques for automated speech assessment and disease-related voice analysis.

Furthermore, Hybrid models combining CNN with recurrent models like LSTM (Long Short-Term Memory), GRU (Gated Recurrent Unit) have been explored to explore spectral-temporal information for the disease classification [2] [4] [9]. Meanwhile, self-supervised representation learning has recently been introduced for automatic speech understanding, which extracts high level semantic representation from unlabeled audio and often utilizes pre-trained speech models to transfer knowledge to downstream tasks for identifying specific vocal conditions such as voice disorders [10]. Compared to these different types of deep models and approaches for automated speech understanding and voice analysis for detection of diseases, we also found that most methods aim to improve classification performance, through complex deep learning models, often sacrificing model interpretability and effective feature learning. Speech analysis to determine PD classification is often challenged by massive data volume due to high-dimensionality, redundancy, and noise.

## 2.1 Research Gap

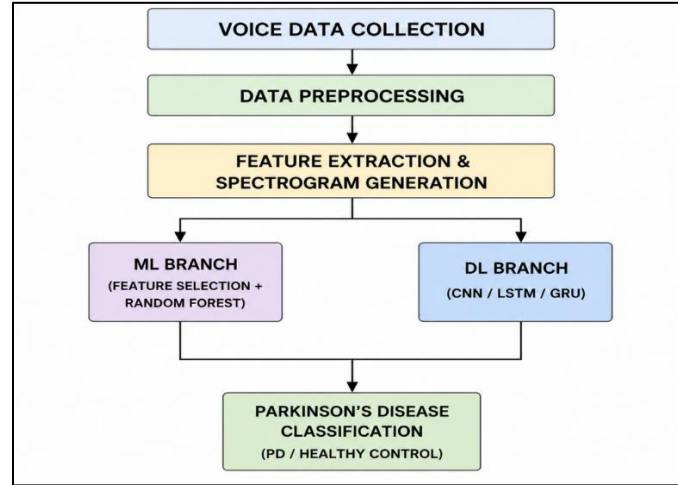
Although machine learning, deep learning, transformer-based, and self-supervised learning paradigms have made substantial advances in speech-based PD detection, a number of limitations still exist in currently available approaches. Existing solutions generally focus on increasing the classification performance by designing ever more complicated models without much consideration of interpretability and analysis of feature importance. The use of high-dimensional acoustic feature vectors may include redundant, correlated, and irrelevant features that cause computational inefficiency, overfitting, and decrease in the generalization ability of models. Besides, traditional methods of feature selection mostly utilize global importance metrics and do not explicitly consider the impact of features on classification in its early decision-making phases when discriminative biomarkers have the most pronounced effect. Therefore, the problem of selecting clinically relevant and computationally efficient speech biomarkers is still far from being solved. There is a need in a feature selection method that would be able to select acoustically relevant features based on their contribution to decision-making in early classification phase and would combine them with the complementary deep spectrograms. Therefore, this study is investigating the integration of Root-Guided Random Forest Feature Selection with deep learning-based classification models to evaluate the effectiveness of discriminative acoustic biomarkers for PD detection. The proposed framework combines interpretable feature selection with deep learning analysis to improve disease discrimination while maintaining model transparency.

## 3. Methodology

The proposed framework includes acoustic feature extraction, discriminative feature selection, and machine learning approaches for automatic detection of Parkinson's disease using speech signals. Figure 1 shows that the framework comprises dataset collection, data preprocessing, feature extraction, Root-Guided feature selection, classification, and performance evaluation processes. In order to describe the Parkinsonian speech features, the framework applies acoustic and spectral speech biomarkers.

Initially, the speech signals collected from public Parkinson's disease datasets are preprocessed through the application of noise removal, voice activity detection, amplitude scaling, and signal normalization. Secondly, the acoustic and spectral features are extracted to encode the phonatory, articulatory, and prosodic problems caused by PD. Thirdly, in order to

eliminate redundant features and select discriminating vocal biomarkers, the Root-Guided Random Forest Feature Selection approach is applied. Finally, the selected optimal feature set is used for classification and prediction tasks.



**Figure 1.** Proposed Workflow of the Root-Guided Parkinson's Disease Detection Framework

### 3.1 Dataset Description

The experimental evaluation was conducted using two publicly available Parkinson's disease speech datasets, namely the PC-GITA dataset [11] and the Parkinson Speech Dataset [12] with Multiple Types of the Sound Recordings. These datasets contain voice recordings collected from PD patients and healthy control (HC) subjects and are widely utilized for speech-based neurological disorder analysis.

The speech recordings include sustained vowel phonation, reading passages, and diadochokinetic (DDK) speech exercises, which are commonly used to capture speech impairments associated with PD. The data sets were pre-processed through a standardized procedure before extracting features and classifying them. The specific details of the data set, pre-processing parameters, feature extraction and selection parameters used in this study are presented in Table 1 below.

**Table 1.** Dataset Characteristics and Experimental Parameters

Parameter	Value
Total Voice Recordings	195
Initial Feature Count	107
Selected Features	Top 75 Features

Feature Selection Method	Root-Guided Random Forest
Random Forest Trees	100
CNN Epochs	50
Learning Rate	0.001
Training–Testing Split	80:20

### 3.2 Data Preprocessing

The speech data recorded from the Parkinson's disease dataset underwent a standard pre-processing stage in order to increase the quality of the signal and ensure uniformity due to differences in the recording environment. As the speech recordings were from various people with different recording intensities, background noises, and lengths, it was necessary to pre-process the data prior to extracting the features.

All the speech recordings were first resampled at a 16 kHz sampling rate to provide uniformity. Next, the spectral noise reduction methods were utilized in order to reduce external interference to avoid any unwanted noise effects when processing features. Then Voice Activity Detection was utilized to detect and remove silence periods and non-speech parts.

Amplitude normalization was done after removing the silent portions in order to counteract the effects of variation brought about by the sensitivity of the microphone and the volume at which the speaker was speaking. The speech signals were then broken down into short frames using Hamming Windowing which eliminates the effect of spectral leakage. This ensured that the representation of the speech signal was robust before feature extraction and Root Guided Feature selection.

### 3.3 Feature Extraction

Feature extraction can help recognize the unique speech characteristics of PD. In particular, given that the speech of Parkinsonian patients is usually characterized by low vocal intensity, impaired articulation, phonatory instability, monotonicity of pitch, and abnormal prosody, an extensive set of acoustic and spectral features was extracted for all speech recordings. Preprocessed speech data were analyzed using both time-domain and frequency-domain methods. Mel-frequency cepstral coefficients (MFCCs) were extracted to describe the relevant speech spectrum, whereas delta and delta-delta MFCCs were calculated to model dynamic speech parameters and articulation. For the analysis of the process of speech

production, measures of voice perturbation, such as jitter and shimmer, were extracted to describe irregular vibration of vocal folds.

The harmonics-to-noise ratio (HNR) was also calculated to characterize voice quality and noise components of speech. Additionally, formants (F1 – F4) and pitch-related features were selected to model articulatory process and vocal tract resonances and prosodic abnormalities typical of PD patients. Spectral and prosodic characteristics were also represented through the extraction of chroma features, spectral contrast, spectral centroid, bandwidth, roll-off, zero-crossing rate (ZCR), short-time energy, and temporal speech features.

Overall, the feature extraction process generated a total of 107 acoustic and spectral attributes for each speech recording. These features formed into the input to the proposed Root-Guided Random Forest Feature Selection (RGRFFS) framework, which identified the most informative vocal biomarkers and reduced feature redundancy before classification.

### 3.4 Root Guided Random Forest Feature Selection

Feature selection is a critical aspect of Parkinson’s disease detection through speech since it involves dealing with very high-dimensional feature space. While the process of feature extraction yields 107 acoustic and spectral features, it is important to note that not all the extracted features have an equal contribution towards disease discrimination. In fact, having too many irrelevant, redundant and weakly informative features might increase the complexity of computations and reduce classification accuracy. In order to solve this problem, a feature selection method known as RGRFFS is adopted.

The proposed approach evaluates each candidate feature by enforcing it as the root splitting attribute of every decision tree within the Random Forest ensemble. Since the root node represents the most influential decision point in a tree, features selected at this level possess stronger class-discriminative capability. For the given feature ( $f_i$ ), the Random Forest model is repeatedly trained and validated, and the corresponding classification performance is recorded. The discriminative strength of each feature is quantified using a Root Importance Score (RIS), defined as:

$$RIS_i = \frac{1}{K} \sum_{k=1}^K Acc_{i,k} \quad (1)$$

where,

- $RIS_i$  = Root Importance Score of the ( $i^{th}$ ) feature
- $Acc_{i,k}$  = classification accuracy obtained during the ( $K^{th}$ ) validation iteration
- $K$  = total number of validation iterations

A higher RIS value indicates the corresponding feature contributes more effectively to the separation of Parkinson's disease and healthy control samples. After computing the RIS for all candidate features, the features are ranked in descending order according to their discriminative capability:

$$F_{rank} = Sort(RIS_1, RIS_2, \dots, RIS_n) \quad (2)$$

where,

- $F_{rank}$  = Ranked feature set
- $RIS_i$  = Root Importance Score of the  $i^{th}$  feature
- $n$  = Total number of features

The optimized feature subset is subsequently constructed by selecting the highest-ranked features:

$$F_{opt} = \{f_1, f_2, \dots, f_m\} \quad (3)$$

where,

- $F_{opt}$  = represents the optimized feature subset
- $f_i$  = Selected feature
- $m$  = Number of selected features

Unlike standard techniques in Random Forest for calculating the importance of each feature considering its total contribution to the entire model, the new method analyzes the potential of each feature for doing the classification job starting from the root node level. Given that the root node is the most significant one in the decision tree, the features picked up at this level have more discriminative power when classifying PD speech versus normal speech of

each feature is evaluated independently and assigned a Root Importance Score (RIS) based on its classification performance across multiple validation iterations.

**Table 2.** Representative Features and Corresponding RIS Scores Identified by the Proposed RGRFFS Framework

Rank	Feature	Feature Category	RIS Score
1	MFCC-1	Spectral Feature	0.94
2	MFCC-2	Spectral Feature	0.93
3	Jitter	Voice Perturbation	0.92
4	Shimmer	Voice Perturbation	0.91
5	HNR	Voice Quality	0.90
6	Pitch Mean	Prosodic Feature	0.89
7	Formant F1	Formant Feature	0.88
8	Formant F2	Formant Feature	0.87

Table 2 highlights the representative acoustic biomarkers which were determined using the Root-Guided Random Forest Feature Selection (RGRFFS) model. Those features that relate to spectral representation, voice perturbation, voice quality, articulation and prosody were found to have higher Root Importance Score (RIS), meaning that they have better discrimination power for PD identification. In particular, those acoustic features based on Mel-Frequency Cepstral Coefficients (MFCCs), jitter, shimmer, Harmonic-to-Noise Ratio (HNR), pitch and formants contributed significantly to differentiating between Parkinsonian and healthy speech samples.

### 3.5 Deep-Learning Based Classification

After applying the Root-Guided Random Forest Feature Selection approach, the optimum set of features having the best discriminative properties were selected to perform classification. Three different types of deep learning techniques have been used for learning complex non-linear relations between speech features to enhance discrimination between PD and HC speech recordings. Since Parkinsonian speech involves subtle changes in phonation, articulation, pitch stability, and voice quality, deep learning models offer a good solution for automatically modeling such complex interactions between the features. In order to test the discriminating abilities of the extracted features, three different types of deep learning models,

CNN, LSTM, and GRU, have been applied. While CNN is used for learning local spectral patterns and feature interactions, LSTM and GRU are chosen to model temporal dependencies between speech features.

To address the imbalance between PD and healthy control samples, class weighting was incorporated during model training. The class weight assigned to each category was computed as:

$$W_j = \frac{N_{\text{total}}}{K \times N_j} \quad (4)$$

Where,

- $W_j$  = weight assigned to the  $j^{\text{th}}$  class
- $N_{\text{total}}$  = total number of samples in the dataset
- $K$  = number of classes
- $N_j$  = number of samples belonging to the  $j^{\text{th}}$  class

The trained models produce class labels for each audio clip, thereby allowing discrimination between the speech audio of Parkinson's patients and the normal control group. From the comparative evaluation carried out on various deep learning models, it was found that CNN had the best classification accuracy. Therefore, CNN became the main classifier of the proposed model, whereas LSTM and GRU were maintained as benchmark classifiers to be compared with the CNN in terms of classification performance.

### 3.6 Algorithm

---

#### Algorithm 1: Root-Guided Random Forest Feature Selection for Parkinson's Disease Detection

---

Input:

- Pre-processed speech feature matrix ( $X$ ) containing 107 acoustic and spectral features
- Class labels ( $Y$ )
- Feature set ( $F = \{f_1, f_2, \dots, f_m\}$ ), (107)
- Number of decision trees ( $T$ )

Output:

- Optimized feature subset ( $F_{\text{opt}}$ ) containing the Top 75 discriminative features

Procedure:

1. Acquire and pre-process speech recordings from the Parkinson's disease datasets.
  2. Extract 107 acoustic and spectral features from each recording.
  3. Construct the complete feature set ( $F$ ).
  4. For each feature ( $f_i \in F$ ):
    - a. Force ( $f_i$ ) as the root splitting attribute in all decision trees.
    - b. Construct a Random Forest ensemble containing ( $T$ ) trees.
    - c. Train the ensemble using the training dataset.
    - d. Evaluate classification performance using validation.
    - e. Compute the Root Importance Score (RIS) of ( $f_i$ ).
  5. Store RIS values for all features.
  6. Rank features in descending order according to RIS.
  7. Select the Top 75 features with the highest RIS values.
  8. Construct the optimized feature subset ( $F_{opt}$ ).
  9. Train classification models using ( $F_{opt}$ ).
  10. Generate Parkinson's disease and healthy control predictions.
  11. Compute classification performance metrics.
  12. Return the optimized feature subset and classification results.
- 

The Root-Guided Random Forest Feature Selection algorithm assesses the discriminative ability of each acoustic feature by estimating the contribution of the feature in making root-level decisions. Contrary to the traditional feature importance estimation techniques which sum up the importance of each feature across several tree levels, the proposed algorithm focuses on the role of features at an early stage of classification. Those features which have higher Root Importance Score (RIS) are characterized by higher class separation ability. After assessing the importance of individual features, the top 75 informative biomarkers are selected from a total number of 107 acoustic and spectral features.

#### 4. Results and Discussion

The performance of the proposed Root-Guided framework was analyzed through several experiments performed on speech data taken from people affected by Parkinson's disease as well as healthy subjects. The performance was quantified by employing popular classification measures such as accuracy, precision, recall, and F1 score, which offer an overall evaluation of the efficiency of the model. Additionally, the comparisons with several machine learning and deep learning frameworks were carried out to analyze the impact of the introduced feature selection method and classification approach. The findings obtained provide

information regarding the discriminative capabilities of the chosen features, the effectiveness of the models used, and the potential of the proposed framework for Parkinson's disease identification.

#### 4.1 Classification Performance Analysis

The proposed Root-Guided Random Forest Feature Selection (RGRFFS) framework were achieved an overall classification accuracy of 93.85%, demonstrating its effectiveness in distinguishing PD speech from healthy control speech samples. As summarized in Table 3, the obtained precision, recall, and F1-score values of 0.94 indicate balanced classification performance and reliable disease discrimination.

**Table 3.** Performance Comparison of Machine Learning and Deep Learning Models

Model	Accuracy (%)	Precision	Recall	F1-Score
CNN	92.3	0.92	0.92	0.92
LSTM	86.1	0.86	0.86	0.86
GRU	83.5	0.83	0.84	0.83
Random Forest (107 Features)	89.2	0.89	0.89	0.89
Root-Guided RF (75 Features)	93.8	0.94	0.94	0.94

The comparative results presented in Table 3 demonstrate the effectiveness of combining Root-Guided feature selection with deep learning-based classification. Among the evaluated deep learning models, CNN achieved the highest classification accuracy of 92.3%, outperforming both LSTM and GRU architectures. The superior performance of CNN can be attributed to its ability to effectively capture local acoustic patterns associated with Parkinsonian speech. Additionally, the proposed Root-Guided Random Forest framework achieved an overall accuracy of 93.85%, demonstrating that discriminative feature selection significantly enhances classification performance.

#### 4.2 Feature Selection Effectiveness Analysis

The efficacy of the developed Root-Guided Random Forest Feature Selection (RGRFFS) method was further validated through assessing the discriminatory power of the selected set of acoustic biomarkers. In particular, the feature selection procedure favored speech features reflecting spectral parameters, voice instability measurements, pitch parameters, and vocal tract resonances that were found to be considerably influenced by PD. Such speech

features as MFCCs, jitter, shimmer, HNR, pitch parameters, and formants exhibited good discriminatory performance owing to the ability of detecting phonatory variability, articulatory impairment, and prosodic disturbances characteristic for Parkinsonian speech.

The selection process helped in minimizing feature redundancy while still maintaining information that is clinically significant to differentiate between the diseased and healthy subjects. Since the proposed method considers acoustic features that are highly discriminating, it helps in improving class separability. The performance evaluation obtained from the proposed system, along with feature ranking results, shows that the extracted features capture the full spectrum of Parkinsonian speech.

### 4.3 Performance Validation and ROC-AUC Analysis

The effectiveness of the proposed Root-Guided Random Forest Feature Selection (RGRFFS) framework was further validated using the classification report presented in Fig. 2

	precision	recall	f1-score	support
Healthy	0.94	0.93	0.94	20
Parkinson	0.94	0.95	0.94	19
accuracy	-	-	0.94	39
macro avg	0.94	0.94	0.94	39
weighted avg	0.94	0.94	0.94	39

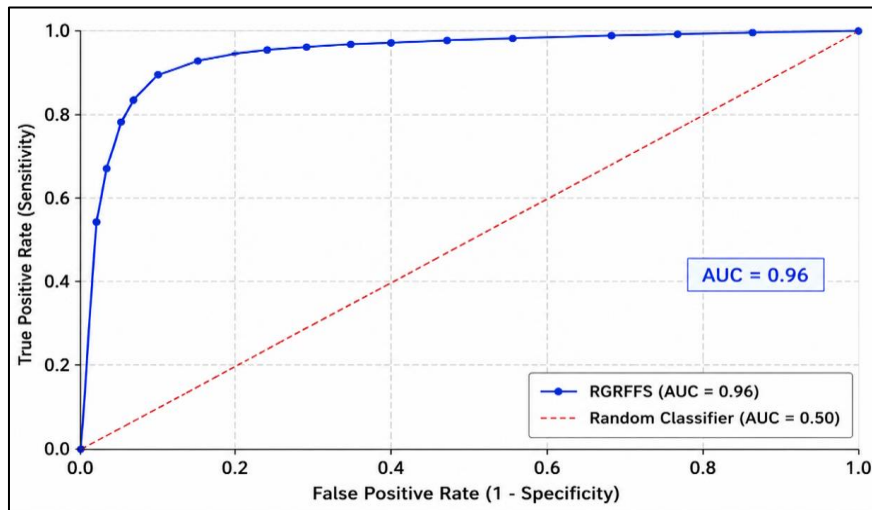
**Figure 2.** Classification Report of the Proposed Root-Guided Random Forest Feature Selection (RGRFFS) Framework

As illustrated in Figure 2, the classifier achieved balanced performance across both PD and healthy control classes. Precision, recall, and F1-score values were consistently maintained at approximately 0.94, indicating stable classification behavior and reliable disease discrimination. For healthy control samples, precision and recall values of 0.94 and 0.93 were obtained, respectively, while PD samples achieved precision and recall values of 0.94 and 0.95. The similarity of these values demonstrates that the framework effectively minimizes classification bias and maintains robust prediction capability across both classes.

**Table 4.** ROC-Based Validation Metrics

Metric	Value
AUC	0.96
Sensitivity	0.95
Specificity	0.95

Table 4 summarizes the diagnostic performance of the proposed framework. An AUC value of 0.96 indicates excellent discrimination between PD and healthy speech samples. The sensitivity and specificity results show that the model can accurately identify both PD cases and healthy individuals. Overall, these results confirm the effectiveness and reliability of the proposed Root-Guided Random Forest Feature Selection framework.



**Figure 3.** ROC Curve of the Proposed Root-Guided Random Forest Feature Selection Framework

The ROC curve shown in Figure 3 provides an additional proof of the efficiency of the suggested technique. Indeed, an Area Under the Curve (AUC) of 0.96 was attained, which demonstrates excellent separation capacity between samples of speech with PD and those with healthy speech. As for the ROC curve, it remains near the upper left part of the coordinate system, which shows high sensitivity and specificity on different decision thresholds. In other words, the selected acoustic biomarkers have excellent discrimination capacity when applied to speech-based PD detection. Obtained AUC value of 0.96 proves excellent discrimination capacity of the suggested approach. An AUC value approaching 1.0 suggests that the classifier can effectively distinguish PD speech from healthy control speech across a wide range of

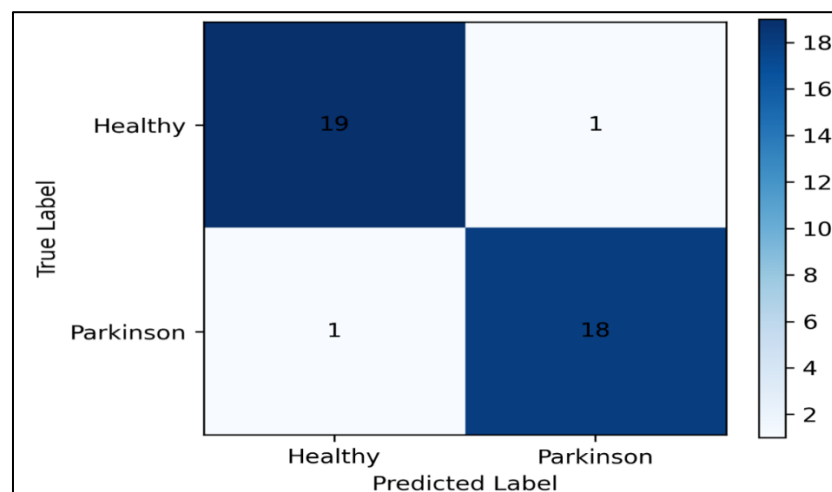
decision thresholds. Furthermore, the ROC curve remains close to the upper-left corner of the coordinate space, indicating a favourable balance between true positive and false positive rates.

The ROC curve analysis is also validated by the outcomes of the classification report and the confusion matrix. The consistent and equal figures of balanced precision, recall, and F1-score values of 0.94 and the small amount of misclassification errors in the confusion matrix show the stability of classification results of the proposed framework in different evaluation metrics. In the clinical aspect, the high AUC score implies that the chosen speech biomarkers have high discriminatory power in determining the features of Parkinsonian speech. Therefore, it can be concluded that the proposed framework can be used as an effective diagnostic tool for PD screening.

The combined results obtained from the classification report and ROC analysis validate the effectiveness of the proposed RGRFFS framework. An overall classification accuracy of 93.85% was achieved. Concurrently, maintaining balanced precision, recall, and F1-score values, demonstrating the suitability of the proposed methodology for reliable and efficient PD detection using speech signals.

#### 4.4 Classification Results Based on the Confusion Matrix

The confusion matrix shown in Fig. 4 provides a detailed assessment of classification performance of the proposed Root-Guided Random Forest Feature Selection (RGRFFS) framework. The matrix illustrates the distribution of true and false classified speech samples for both PD and healthy control categories.



**Figure 4.** Confusion Matrix of the Proposed Root-Guided Random Forest Feature Selection (RGRFFS) Framework

As observed in Figure 4, there have been 19 samples of healthy speech that have been classified correctly as healthy, but there was only one case where healthy speech was wrongly classified as having PD. There have also been 18 samples of PD speech that were correctly identified, but there was only one sample where PD speech was wrongly classified as healthy.

The confusion matrix also shows that the classification performance is balanced across the two classes. This can be seen by the large number of data instances concentrated around the primary diagonal, which is an indication of agreement between the true and the predicted class labels. These findings are consistent with the obtained precision, recall, and F1-score values of 0.94, further validating the reliability of the proposed framework for speech-based PD detection.

## 5. Conclusion and Future Work

The Root-Guided Random Forest Feature Selection (RGRFFS) algorithm for speech-based Parkinson's disease recognition is discussed in this paper. In this method, the feature selection process used acoustic and spectral speech parameters to discover biomarkers linked to Parkinsonian speech. Through ranking the features based on their Root Importance Score (RIS), the most discriminating attributes were selected without loss of classification performance and with reduction of redundant features. The findings suggest that Root-Guided feature selection technique enables an efficient feature selection method for discriminative acoustic biomarkers in PD classification. It was found out that the framework is able to achieve accuracy of 93.85%, with precision, recall, and F1 score equal to 0.94, demonstrating its ability to discriminate between speech of patients with PD and normal speech with high reliability. The research shows the efficiency of the feature selection technique in terms of improvement of classification performance and disease discrimination. Future work may include the use of large and varied speech data sets for enhanced performance and generalization. In addition, deep learning algorithms and multimodal biomarkers may be investigated for use in real-time screening of Parkinson's disease.

## References

- [1] Madusanka, Nuwan, and Byeong-il Lee. "Vocal Biomarkers for Parkinson's Disease Classification Using Audio Spectrogram Transformers." *Journal of Voice* 2024.

- [2] Wodzinski, Marek, Andrzej Skalski, Daria Hemmerling, Juan Rafael Orozco-Arroyave, and Elmar Nöth. "Deep Learning Approach to Parkinson's Disease Detection Using Voice Recordings and Convolutional Neural Network Dedicated to Image Classification." In 2019 41st annual international conference of the IEEE engineering in medicine and biology society (EMBC), IEEE, 2019, 717-720.
- [3] Hossain, Mohammad Amran, and Francesco Amenta. "Machine Learning-Based Classification of Parkinson's Disease Patients Using Speech Biomarkers." *Journal of Parkinson's Disease* 2024, vol. 14, no. 1: 95-109.
- [4] Chen, Wenna, Rongfu Lv, Xiaowei Du, Xiangyu Chen, Hao Wang, Jincan Zhang, and Ganqin Du. "Parkinson's Disease Detection Using Spectrogram-Based Multi-Model Feature Fusion Networks." *Frontiers in Neurology* 2025, vol. 16: 1706317.
- [5] Hernandez, Abner, Eunjung Yeo, Kwanghee Choi, Chin-Jou Li, Zhengjun Yue, Rohan Kumar Das, Jan Ruzs et al. "Adapting Self-Supervised Speech Representations for Cross-lingual Dysarthria Detection in Parkinson's Disease." arXiv preprint arXiv:2603.22225 (2026).
- [6] Klempíř, Ondřej, and Radim Krupička. "Analyzing Wav2vec 1.0 Embeddings for Cross-Database Parkinson's Disease Detection and Speech Features Extraction." *Sensors* 2024, vol. 24, no. 17: 5520.
- [7] Sedigh Malekroodi, Hadi, Nuwan Madusanka, Byeong-il Lee, and Myunggi Yi. "Speech-Based Parkinson's Detection Using Pre-Trained Self-Supervised Automatic Speech Recognition (ASR) Models and Supervised Contrastive Learning." *Bioengineering* 2025, vol. 12, no. 7: 728.
- [8] Skaramagkas, Vasileios, Anastasia Pentari, Zinovia Kefalopoulou, and Manolis Tsiknakis. "Multi-Modal Deep Learning Diagnosis of Parkinson's Disease—A Systematic Review." *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 2023, vol. 31: 2399-2423.
- [9] Shibina, V., and T. M. Thasleema. "A Hybrid Approach to Detecting Parkinson's Disease Using Spectrogram and Deep Learning CNN-LSTM Network." *International Journal of Speech Technology* 2024, vol. 27, no. 3: 657-671.

- [10] Ribas, Dayana, Miguel A. Pastor, Antonio Miguel, David Martínez, Alfonso Ortega, and Eduardo Lleida. "Automatic Voice Disorder Detection Using Self-Supervised Representations." *IEEE Access* 2023, vol. 11: 14915-14927.
- [11] Orozco-Arroyave, Juan Rafael, Julián David Arias-Londoño, Jesús Francisco Vargas-Bonilla, María Claudia Gonzalez-Rátiva, and Elmar Nöth. "New Spanish Speech Corpus Database for the Analysis of People Suffering from Parkinson's Disease." In *Lrec* 2014, vol. 14: 342-347.
- [12] Little, M. "Parkinsons" [Dataset]. UCI Machine Learning, 2007, [https://archive.ics.uci.edu/dataset/174/parkinsons?utm\\_source](https://archive.ics.uci.edu/dataset/174/parkinsons?utm_source)