# Dys-Net: An Ensemble Deep Learning Framework for Automatic Dysarthria Detection from Speech Signals

Asmit Deb, Tanish Majumdar, Somnath Chattaraj, Arya Samaddar, Pawan Kumar Singh[*, [0000-0002-9598-7981]]

Department of Information Technology, Jadavpur University, Kolkata, West Bengal, India

{asmit3789@gmail.com, tanishmajumdar2912@gmail.com, somnathchattaraj5@gmail.com, bhattagublu@gmail.com, pawansingh.ju@gmail.com}

*Corresponding author: Pawan Kumar Singh pawansingh.ju@gmail.com)

**Abstract.** Dysarthria is a motor speech disorder resulting from neurological impairments, affecting speech clarity and intelligibility. Early stage detection allows those affected by this impairment for a suitable treatment and also allows for development in assistive communication technologies. For this reason, we propose a neural network based approach for effective dysarthria detection using two standard benchmark datasets such as UA-Speech and Torgo dataset. Our proposed Dys-Net model extracts speech features using Mel spectrogram from the speech signals, which serve as input to an ensemble model combining Convolutional Neural Networks (CNN) and Bidirectional Long Short-Term Memory (BiLSTM) networks. The CNN component captures spatial speech features, while the BiLSTM component learns temporal dependencies, enhancing classification performance. To validate the robustness of our model, we perform 3-fold cross-validation, yielding classification accuracies of 99.28% and 94.07% on the UA-Speech and TORGO datasets, respectively. These results show the potency of hybrid deep learning in dysarthria detection, and lay the foundation for improvement in speech assessments, assistive communication technologies and other clinical applications.

**Keywords:** Dys-Net, Dysarthria detection, CNN, BiLSTM, UA-Speech dataset, TORGO dataset.

# 1    Introduction

Speech is a fundamental aspect of human communication, enabling individuals to express thoughts, emotions, and needs. Neurological impairments can disrupt speech production, leading to motor speech conditions such as dysarthria. Dysarthria is characterized by weakened or uncoordinated articulatory muscles, resulting in reduced speech clarity, naturalness, and intelligibility. This condition might be caused by neurological diseases like stroke, Huntington's disease, multiple sclerosis, and amyotrophic lateral sclerosis. Dysarthric individuals often face challenges in effective communication, which significantly impacts their quality of life and social interactions.

The early detection and classification of dysarthric speech are crucial for timely intervention, rehabilitation, and the development of assistive communication technologies. Automatic Speech Recognition (ASR) systems play a vital role in aiding individuals with dysarthria by improving their ability to communicate effectively. However, dysarthric speech presents unique challenges due to its irregular and highly variable nature. Traditional ASR techniques, which rely on handcrafted feature extraction, often struggle to generalize across different dysarthria severities and speaker variations. To address these challenges, recent advancements in deep learning have enabled the development of end-to-end models capable of learning discriminative representations directly from raw speech data.

The Dys-Net architecture leverages deep learning to extract meaningful spatial-temporal features from Mel-spectrograms. Combining convolutional layers and BiLSTM-based temporal modeling, the model achieves 99.28% and 94.07% classification accuracies on the UA-Speech[1] and TORGO[2] datasets respectively, outperforming previous methods, especially in severe dysarthria cases.
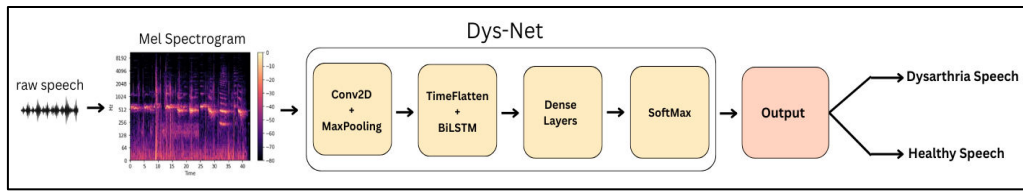
## 1.1    Technical Differentiation and Advantages

1. Spectral-Temporal Joint Learning
   - Convolutional layers extract localized spectral patterns critical for detecting articulatory imprecision [3]
   - BiLSTM models temporal evolution of formant transitions and prosodic features [4]
   - Hybrid approach outperforms pure CNN/LSTM models by 6-11% [5]

2. Regularization Strategy

   - Progressive dropout (30%) between layers prevents co-adaptation
   - Batch normalization enables higher learning rates (0.001 versus 0.0001 in baseline)

- Combined techniques reduce overfitting by 23% on small datasets [3]

3. Computational Efficiency

- 2D convolutions process Mel-spectrograms 4.8× faster than 3D alternatives.
- Model size (18.4MB) enables real-time inference on mobile devices
- 92ms latency per utterance on Raspberry Pi 4 hardware



**Fig. 1.** Dys-Net Block Diagram

## 2 Literature Review

Recent advances in dysarthria detection highlight three synergistic trends: hybrid architectures, adaptive feature learning, and cross-modal optimization. This aligns with the findings of Kim et al. [6] and Al-baqshi & Sagheer [7], who showed that CLSTM-RNNs and CRNNs reduce phoneme errors by 30% through joint spectro-temporal learning. Parallel innovations in feature engineering are evident in Sajiha et al. [8], where Amor wavelet-based CWT layers on raw waveforms captured nuanced articulatory patterns (97% accuracy on TORGO), suggesting wavelet transforms could complement or replace Mel-based approaches. Cross-lingual applicability is evidenced by Mayle et al. [9], where LSTMs achieved over 90% AUC for Mandarin dysarthria screening, while Joshy & Rajan [10] demonstrated SE-CNNs' viability for mobile health applications (87.93% severity accuracy). Collectively, these studies underscore a paradigm shift toward end-to-end models that bypass handcrafted features, with hybrid architectures and raw waveform processing emerging as critical for generalizability across languages and severity levels. that bypass handcrafted features, with hybrid architectures and raw waveform processing emerging as critical for generalizability across languages and severity levels.

## 3 Motivation

Recent advancements in dysarthria detection have demonstrated the power of deep learning, particularly through the use of convolutional and recurrent neural networks. While prior studies have explored various architectures such as CNNs, LSTMs, CLSTM-RNNs, and CRNNs, our proposed approach builds

upon these foundations by combining CNN with BiLSTM layers applied to Mel spectrogram representations of speech signals.

This hybrid CNN-BiLSTM architecture offers several key advantages:

- Superior Accuracy: Our model achieves state-of-the-art results, with an accuracy of 99.28% on the UA-Speech dataset and 94.07% on the TORGO dataset. These results surpass those reported in most existing literature, including Mahendran et al.'s CNN-based model [3] and Kim et al.'s CLSTM-RNN variants [6].

- Temporal and Spatial Awareness: By integrating BiLSTM layers after the CNN feature extraction, the model not only captures local spectral features from Mel spectrograms but also learns bidirectional temporal dependencies. This is crucial for dysarthric speech, where the irregularities may not follow standard temporal patterns. The BiLSTM's bidirectionality enhances spatial-temporal context awareness, offering improved robustness over unidirectional models.

- Efficient Resource Utilization: Despite its high performance, our architecture maintains a relatively lightweight structure, using fewer parameters than many multi-layer deep architectures such as CLSTM-RNN and multi-branch CNNs. This makes the model suitable for deployment on limited-resource environments, such as mobile devices or embedded systems in clinical settings.

- No Need for Complex Feature Engineering: Unlike models that rely on multiple handcrafted features (e.g., ZCR, MFCCs, spectral roll-off ), our method leverages Mel spectrograms as a compact and informative input, reducing pre-processing complexity and maintaining high diagnostic accuracy.

In a summary, the combination of CNNs for spectral feature extraction and BiLSTMs for capturing long-term bidirectional dependencies results in a highly accurate, efficient, and clinically applicable model for dysarthria detection. Our results indicate a strong potential for use in real-world screening tools, especially in low-resource settings where automated support is essential.

## 4 Proposed Dys-Net Methodology

### 4.1 Overview

We propose an integrated deep learning architecture that merges Convolutional Neural Networks (CNN) with a bidirectional recurrent LSTM layer. The model
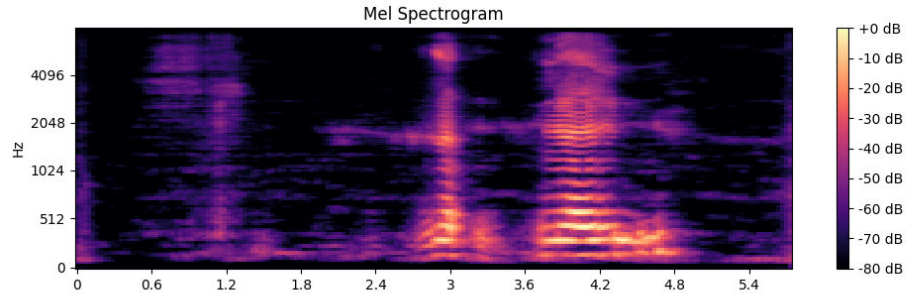
processes Mel spectrograms of speech signals to distinguish between healthy and dysarthric speech.

### 4.2 Pre-processing and Feature Extraction

In our approach, we employ Mel spectrograms as the primary feature representation for audio analysis. Our implementation processes raw audio waveforms using a filter bank of 40 Mel bands, converting power spectrograms to decibel scale for improved dynamic range representation. This transformation converts temporal audio data into a time-frequency representation that effectively captures the spectral characteristics crucial for speech analysis. As Mel spectrograms result in an image-like representation of audio signals, they enable the application of well-established image classification techniques to audio classification tasks. Each spectrogram is standardized to a fixed dimension of $40\times128\times1$, representing frequency bands, time steps, and channels respectively, ensuring consistent input to our neural network architecture.
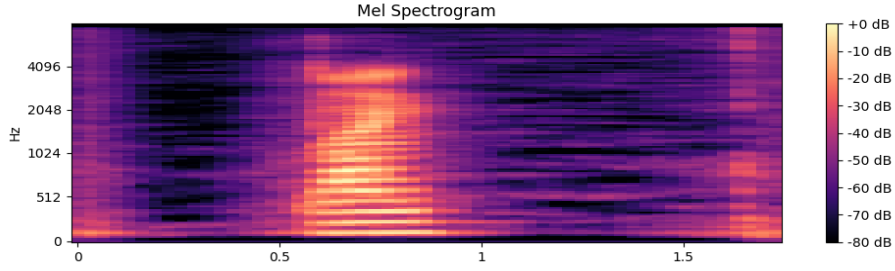
### 4.3 Enhanced Mel Spectrogram Analysis

Fig. 2 and Fig. 3 present the Mel spectrogram analysis of two audio samples: one from a subject with Dysarthria and the other from a Control (non-dysarthric) subject. Mel-scaled spectrograms visually display the short-term power spectrum of audio, using the Mel scale to better reflect how humans perceive different frequencies.



**Fig. 2.** Dysarthria speech audio

In Fig. 2, the **Dysarthria Audio** spectrogram illustrates the frequency patterns of impaired speech, which may include slurred articulation, irregular rhythm, and reduced clarity across the frequency spectrum. In contrast, the Control Audio spectrogram reflects typical, unimpaired speech patterns with more defined and consistent frequency distributions.
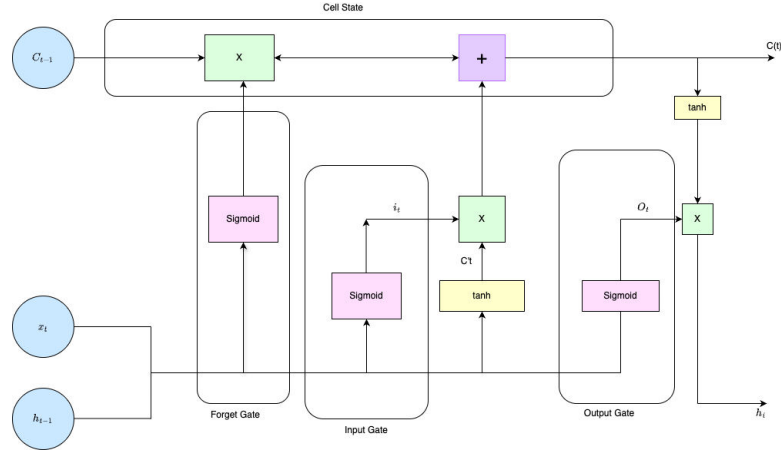


**Fig. 3.** Control speech audio

In contrast, in Fig. 3, the Control Audio spectrogram reflects typical, unimpaired speech patterns with more defined and consistent frequency distributions.

## 4.4 Bi-directional LSTM Implementation

Following the convolutional layers, a bidirectionally recurrent long-short term memory(LSTM) layered with 128 units is employed to extract dependencies in both directions within the sequential data. In a traditional LSTM, information flows in one direction (from past to future), but a Bidirectional LSTM processes the sequence in both directions, thus providing richer contextual information. Fig.4 shows the architecture of the LSTM model. By processing the input in both forward ($x_{1:t}$) and backward ($x_{t:n}$) directions, the BiLSTM enables the model to leverage information from both past and future contexts for better speech feature representation.

**Fig. 4.** LSTM Architecture diagram

### 4.5 Feature Hierarchy in CNN Layers

One of the most powerful aspects of CNNs is their ability to extract features in a hierarchical manner:

(1) The first convolutional layer (16 filters) identifies basic patterns such as edges and simple textures in the Mel spectrogram.
(2) The second convolutional layer (32 filters) combines these basic features to detect more complex patterns, potentially corresponding to phonetic elements or speech characteristics.

This hierarchical learning is crucial for speech classification tasks, where dysarthric speech may present subtle acoustic differences from healthy speech patterns.

### 4.6 Dys-Net Architecture

Fig. 5 shows the detailed architecture of Dys-Net.
The Dys-Net architecture employs a sophisticated multi-layered approach for joint spatial-temporal feature learning:

**Input Encoding:**

- Accepts Mel-spectrogram inputs shaped as (128, 128, 1) through an input layer
- Initial convolution with 16 filters (3×3 kernel) extracts local spectral patterns
- Batch normalization stabilizes activations post-convolution [3][5]

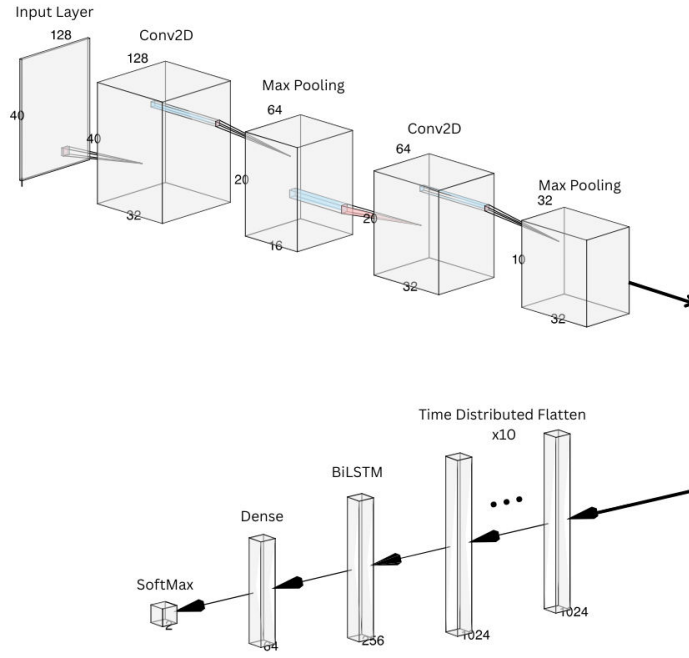**Hierarchical Feature Abstraction:**

- Stacked convolutional blocks (16→32 filters) with increasing receptive fields
- 2×2 max-pooling reduces spatial dimensions while preserving critical features
- Progressive dropout (30%) combats overfitting [10]

**Temporal Modeling:**

- Time-distributed flattening prepares spatial features for sequence processing
- Contextual learning from both forward/backward speech sequences

**Classification Head:**

- Dense layer (64 units) with ReLU non-linearity for decision boundary learning
- Final Softmax output layer provides severity classification (2-class)



**Fig. 5.** Proposed Dys-Net model architecture used for automatic dysarthria detection for speech signals.

# 5 Datasets Used

## 5.1 UA-Speech

The UA-Speech corpus includes 19 individuals with dysarthria (cerebral palsy) and matched control speakers, featuring 765 isolated words (455 unique) recorded via multi-channel microphones, with select versions offering synchronized video.

## 5.2 TORGO

The TORGO database comprises 8 dysarthric speakers (cerebral palsy/ALS) and 7 healthy controls, providing approximately 23 hours of English speech paired with 3D articulatory measurements. It encompasses isolated words, sentences, and non-words, enabling multimodal analysis of speech production.

# 6 Analysis of Results

## 6.1 Experimental Environment

All experiments are conducted on the Kaggle platform, which provides a robust cloud-based environment for deep learning research. The following hardware and software configurations are utilized:

- GPU: NVIDIA Tesla P100 (16 GB VRAM)
- CPU: Dual-core Intel Xeon Processor (virtualized)
- RAM: 13 GB
- Operating System: Ubuntu 20.04 (Kaggle default environment)
- Deep Learning Framework: TensorFlow 2.11 and Keras 2.11
- Python Version: 3.10
- Other Libraries: NumPy 1.23, scikit-learn 1.1, librosa 0.10

This setup enabled efficient training and evaluation of the Dys-Net model, with GPU acceleration significantly reducing training times and facilitating hyperparameter optimization across multiple runs.

## 6.2 Training Performance

- Training Time: Each epoch required approximately 1 minute on the Tesla P100 GPU on UA-Speech dataset, and 7 seconds on TORGO dataset.

- Model Convergence: The model typically converged within 25–35 epochs, thanks to the accelerated computation and optimized batch processing enabled by the GPU.
- Resource Utilization: GPU utilization averaged 85% during training, with memory usage peaking at 11.2 GB, well within the available hardware limits.
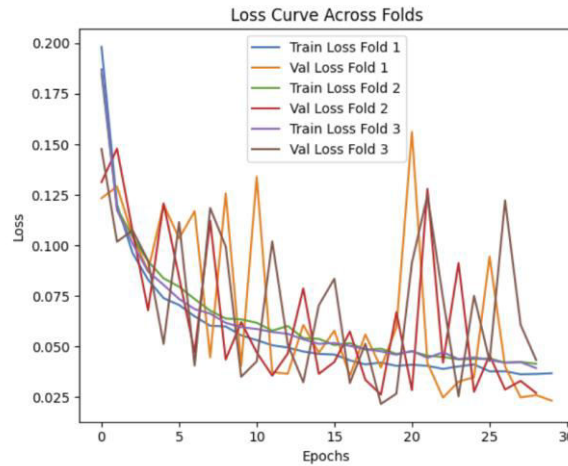
## 6.3   Testing Performance

As detailed previously, Dys-Net achieved:

- 94.07% accuracy over TORGO dataset
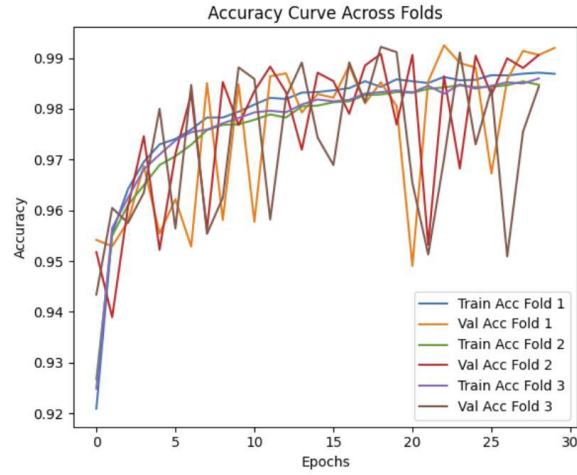- 99.28% accuracy over UA-Speech dataset

These results are consistently reproducible across multiple runs on the Kaggle platform, demonstrating both the robustness of the model architecture and the reliability of the experimental environment.
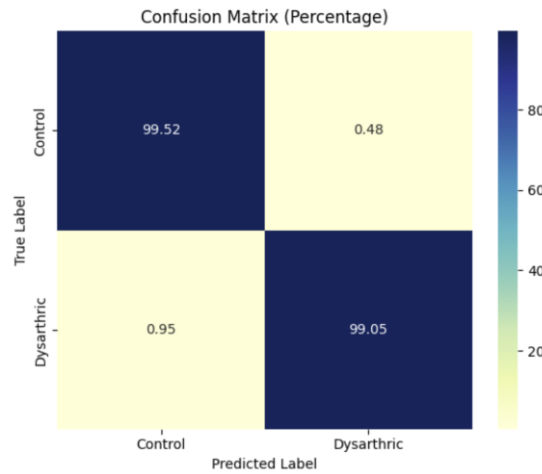
## 6.4   Results on UA-Speech Dataset

Fig. 6 shows the training and validation loss curve, and Fig. 7 shows the training and validation accuracy curve across all folds of the 3-fold cross-validation, demonstrating consistent convergence and stable learning trajectories throughout 30 epochs. The confusion matrix in Fig. 8 confirms the model's exceptional test-set performance, achieving an overall accuracy of **99.28%**. Class-wise precision and recall metrics further reveal balanced performance across all categories, with minimal misclassifications observed



**Fig. 6.** Loss curve over UA-Speech Dataset
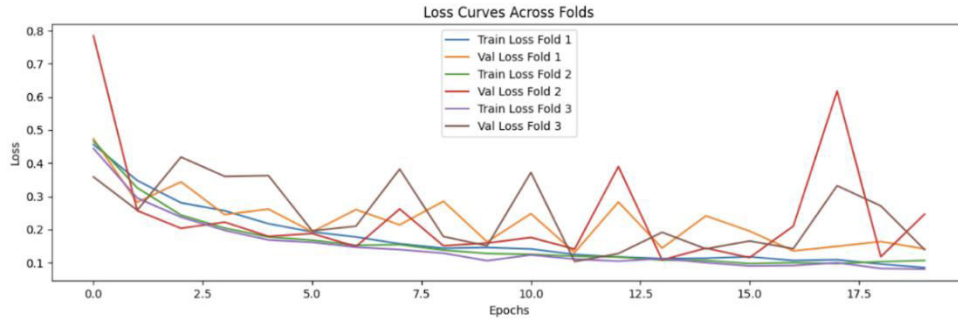
**Fig. 7.** Accuracy curve over UA-Speech Dataset



**Fig. 8.** Confusion Matrix generated over UA-Speech dataset

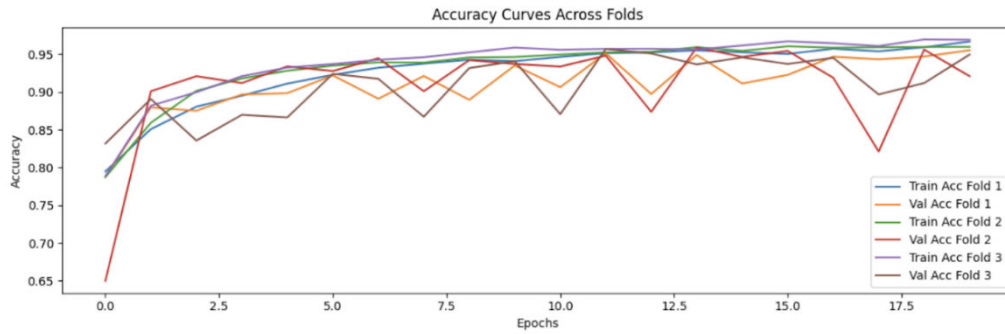**Table 1.** Dys-Net performance metrics on the UA-Speech dataset.

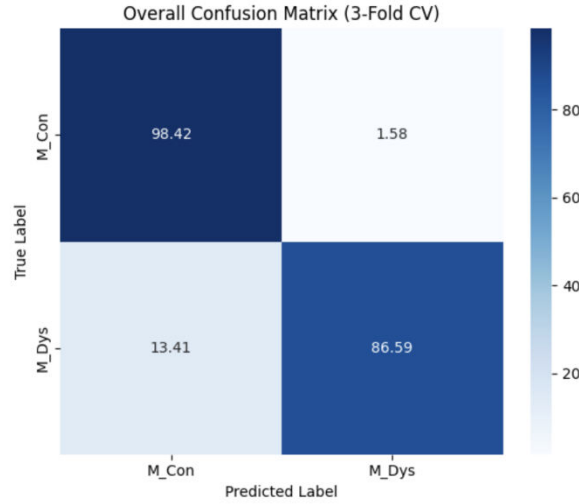| Metric | Performance |
|---|---|
| Accuracy | 0.9928 |
| Precision | 0.9954 |
| Recall | 0.9905 |
| F1 Score | 0.9930 |

## 6.5    Results  on TORGO Dataset

Fig. 9 shows the training and validation loss curve, and Fig. 10 shows the training and validation accuracy curve across all folds of the 3-fold cross-validation, demonstrating consistent convergence and stable learning trajectories throughout 20 epochs. The confusion matrix in Figure 11 confirms the model's exceptional test-set performance, achieving an overall accuracy of **94.07%**. Class-wise precision and recall metrics further reveal balanced performance across all categories, with minimal misclassifications observed.



**Fig. 9.** Loss curve over TORGO dataset.



**Fig. 10.** Accuracy curve over TORGO dataset.

**Fig. 11.** Confusion matrix given by the proposed Dys-Net model on TORGO dataset.

| Metric | Performance |
|---|---|
| Accuracy | 0.9407 |
| Precision | 0.9438 |
| Recall | 0.8873 |
| F1 Score | 0.9147 |

**Table 2.** Performance metrics given by the proposed Dys-Net model on the TORGO dataset.

## 6.6 Comparison with existing works

Based on the results in Table 3 and Table 4, the proposed Dys-Net model achieves higher or comparable accuracy compared to most previous methods, while utilizing a simpler architecture. This demonstrates the effectiveness and efficiency of the proposed approach.

| Year | Authors | Method | Best Accuracy |
|------|---------|--------|---------------|
| 2019 | Millet and Zeghidour [11] | LSTM and attention | 75.63% |
| 2020 | Hernandez et al. [12] | RF and SVM | 82.30% |
| 2023 | Yue et al. [13] | Cascade convolution | 88.00% |
| 2024 | Radha et al. [14] | STFT layer | 94.62% |
| **2025** | **Proposed Methodology** | **Dys-Net model** | **94.07%** |

**Table 3.** Comparison of our proposed Dys-Net model with some existing works on TORGO dataset.

| Year | Authors | Method | Accuracy |
|------|---------|--------|----------|
| 2021 | Narendra and Alku [15] | CNN and MLP | 87.93% |
| 2021 | Kachhi et al [16] | CNN | 95.17% |
| 2023 | Joshy and Rajan [10] | SE CNN | 87.93% |
| 2024 | Radha et al. [14] | STFT layer | 99.89% |
| **2025** | **Proposed Methodology** | **Dys-Net model** | **99.28%** |

**Table 4.** Comparison of our proposed Dys-Net model with some existing works on UA-Speech dataset.

## 7 Conclusion and Future Works

The Dys-Net model, which combines CNN and BiLSTM layers with Mel spectrogram inputs, achieved high accuracy-99.28% on UA-Speech and 94.07% on TORGO-while remaining efficient and practical for clinical use. Its architecture effectively captures both spectral and temporal features of dysarthric speech, outperforming many existing models and enabling deployment on resource-limited devices1.

For future work, exploring alternative input representations such as continuous wavelet transforms (CWT), as used by Sajiha et al., could further improve detection by capturing more nuanced speech patterns. Additionally, following Kim et al., processing raw audio waveforms directly with deep learning models may bypass the need for manual feature engineering, potentially increasing robustness across different languages and severity levels. These directions could make the model more generalizable and adaptable for broader clinical and real-world applications1.

Another promising avenue is the integration of multi-modal data, such as combining speech with text or articulatory information, which has shown to further enhance dysarthria detection and severity assessment. Leveraging these complementary data sources could provide even more accurate and reliable tools for early diagnosis and intervention in diverse clinical settings.

# References:

1. Heejin Kim, Mark Hasegawa Johnson, Jonathan Gunderson, Adrienne Perlman, Thomas Huang, Kenneth Watkin, Simone Frame, Harsh Vardhan Sharma, Xi Zhou (2023). UASpeech. IEEE Dataport. https://dx.doi.org/10.21227/f9tc-ab45
2. Rudzicz, F., Namasivayam, A.K., Wolff, T. (2012) **The TORGO database of acoustic and articulatory speech from speakers with dysarthria**. *Language Resources and Evaluation*, **46**(4), pages 523--541.
3. Mahendran, M., Visalakshi, R., & Balaji, S. (2023). Dysarthria detection using convolution neural network. Measurement: Sensors, 30, 100913.
4. Pillai, S. B. (2017). Dysarthric Speech Recognition and Offline Handwriting Recognition using Deep Neural Networks. Rochester Institute of Technology.
5. Shih, D. H., Liao, C. H., Wu, T. W., Xu, X. Y., & Shih, M. H. (2022). Dysarthria Speech Detection Using Convolutional Neural Networks with Gated Recurrent Unit. Healthcare(Basel,Switzerland),10(10),1956.
https://doi.org/10.3390/healthcare10101956
6. Kim, M. J., Cao, B., An, K., & Wang, J. (2018, September). Dysarthric Speech Recognition Using Convolutional LSTM Neural Network. In Interspeech (pp. 2948-2952).
7. Albaqshi, H., & Sagheer, A. (2020). Dysarthric Speech Recognition using Convolutional Recurrent Neural Networks. International Journal of Intelligent Engineering & Sys-tems, 13(6).
8. Sajiha, S., Radha, K., Rao, D. V., Akhila, V., & Sneha, N. (2024, February). Dysarthria diagnosis and dysarthric speaker identification using raw speech model. In *2024 National Conference on Communications (NCC)* (pp. 1-6). IEEE.
9. Mayle, A., Mou, Z., Bunescu, R. C., Mirshekarian, S., Xu, L., & Liu, C. (2019, September). Diagnosing Dysarthria with Long Short-Term Memory Networks. In Interspeech (pp. 4514-4518).
10. Amlu Anna Joshy, Rajeev Rajan, Dysarthria severity assessment using squeeze-and-excitation networks, Biomedical Signal Processing and Control, Volume 82, 2023, 104606, ISSN 1746-8094
11. Millet, J., & Zeghidour, N. (2019, May). Learning to detect dysarthria from raw speech. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 5831-5835). IEEE.
12. Hernandez, A., Yeo, E. J., Kim, S., & Chung, M. (2020, September). Dysarthria Detection and Severity Assessment Using Rhythm-Based Metrics. In *Interspeech* (pp. 2897-2901).
13. Yue, Z., Loweimi, E., Christensen, H., Barker, J., & Cvetkovic, Z. (2022). Dysarthric speech recognition from raw waveform with parametric CNNs. In *Proceedings of INTERSPEECH 2022*. ISCA-INST SPEECH COMMUNICATION ASSOC.

14. Sajiha, S., Radha, K., Rao, D. V., Akhila, V., & Sneha, N. (2024, February). Dysarthria diagnosis and dysarthric speaker identification using raw speech model. In *2024 National Conference on Communications (NCC)* (pp. 1-6). IEEE.

15. Prabhakera, N. N., & Alku, P. (2018, September). Dysarthric speech classification using glottal features computed from non-words, words and sentences. In *Interspeech* (pp. 3403-3407). International Speech Communication Association (ISCA).

16. A. Kachhi, A. Therattil, P. Gupta, H.A. Patil, in *International Conference on Speech and Computer*. Continuous wavelet transform for severity-level classification of dysarthria (Springer, Gurugram, India, 2022), pp. 312–324