

## **Dysarthric Speech Detection Using CNN & Hybrid Models**

### **Team - 13**

<b>Team Members</b>	<b>Registration Number</b>
SAKILAM ABHIMAN	CH.EN.U4AIE20055
SARTHAK YADAV	CH.EN.U4AIE20058
SHAIK HUZAIFA FAZIL	CH.EN.U4AIE20060
B SIVA JYOTHI NATHA REDDY	CH.EN.U4AIE20063
VEMIREDDY ANVITHA	CH.EN.U4AIE20076

## **Abstract**

Dysarthria is an inability of a person to speak properly because the muscles used for speech, including the tongue, lips, vocal cords, and diaphragm are affected. It can result in difficulty pronouncing words, slurred speech, a slow or fast rate of speech, monotone or unmodulated voice, and problems with breathing while speaking. Dysarthria can be caused by a variety of factors, including neurological conditions such as stroke, multiple sclerosis, or cerebral palsy, as well as traumatic brain injury or certain medications.

Developing a model for detecting dysarthric individuals is an important task as it can aid in early diagnosis, treatment planning, and assistive technology development. In this project, machine learning models such as CNN (Convolutional Neural Network), CNN-LSTM (Convolutional Neural Network-Long Short-Term Memory), and CNN-GRU (Convolutional Neural Network-Gated Recurrent Unit) will be developed to detect dysarthric speech. Of these we have the accuracies as 96.5,97,97.5. We can detect the impaired speech based on acoustic features such as pitch, intensity, and formants. The model will be trained on a dataset of speech recordings from individuals with dysarthria and control participants. The performance of the model will be evaluated using metrics such as accuracy, precision, recall, and F1-score. The ultimate goal of this project is to develop a reliable and accurate tool for detecting dysarthric individuals that can be used in clinical settings and assistive technology development.

Further, we can use this methodology in developing an ASR, which would become an essential technology for many applications such as dictation, transcription, voice assistants, and speech-enabled systems. However, the performance of ASR is challenged when it comes to individuals with speech impairments such as dysarthria, a motor speech disorder that affects the articulation, fluency, and intelligibility of speech. This report includes the datasets used by various papers and methodologies. And it gives what the future holds with respect to dysarthric disease.

## Introduction

Dysarthria is a speech impairment that impairs a person's capacity for clear, understandable speaking. It is brought on by injury to the brain or the nerves that manage the vocal cords, tongue, lips, and other speech-related muscles. A person's capacity to interact successfully with others can be significantly impacted by dysarthria, which can result in frustration, social isolation, and a lower quality of life. Dysarthria can also make it challenging for sufferers to take part in regular activities including job, school, and social gatherings.

Dysarthria is a serious health problem that affects people of all ages and backgrounds since it can be brought on by a wide range of underlying diseases, including cerebral palsy, stroke, brain injury, Parkinson's disease, and multiple sclerosis. There is presently no treatment for dysarthria, while it may be treated with medication, speech therapy, and other supportive measures. Dysarthria can have a wide range of causes and can affect persons of all ages. The below figure shows some typical causes of dysarthria:



**Fig - 01**

Automatic Speech Recognition (ASR) systems have been widely used to enable people with disabilities to interact with technology and the environment. One group of people who benefit from ASR are those with dysarthria, a motor speech disorder caused by damage to the brain or nervous system. In recent years, there have been significant advancements in ASR technology specifically designed for dysarthric people.

One of the recent developments in ASR for dysarthric people is the use of deep learning models, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), to improve the accuracy of recognition. These models can learn complex patterns in speech data and adapt to individual variations in dysarthric speech. Several datasets have been used to train and evaluate ASR systems for dysarthric people. The Dysarthric Speech Corpus (DiSC), a publicly available dataset, contains speech recordings of dysarthric speakers with different degrees of severity. Another dataset is the Scottish Dysarthria Corpus, which includes speech samples from people with a range of neurological conditions causing dysarthria.

In this report simple CNN model, and CNN-LSTM model, CNN-GRU model have used for training and testing on a simple dysarthric speech dataset called Torgo dataset. We attain 3 results, therefore we compare all the 3 models and perform comparative analysis by using different evaluation metrics such as F1\_score, confusion matrix, accuracy.

## Literature Survey

A complete survey has been done on the recent developments on ASR for Dysarthric People. I have used “Google Scholar” for papers, I used different keywords for search and I identified some thousands of papers based on my searches. Later I performed screening and eliminated as many as papers which are not relevant to my topic. Finally, I ended up 40 papers and have chosen 25 of them which are quite relevant to my topic, for the literature review. Those 25 papers are from 2013 to 2023. These papers have clearly mentioned the methods, results, limitations, drawbacks, reasons related to ASR for Dysarthric People. I have considered each paper and have gone through the entire paper, primarily focusing on abstract, results and conclusions. I have encountered similar methods, methodologies in some papers. Also the reasons for their unexpected results are also similar in some papers. Some papers focus on how an ASR is developed, some focus on how it functions, the rest specifies the issues, difficulty involved in making an ASR. Likewise, after going through several kinds of papers, I have made this review which focuses on developments on ASR for Dysarthric People since 2013.

One of the common techniques for ASR for dysarthric individuals is speaker adaptation. Speaker adaptation is the process of adapting the ASR system to the unique speech characteristics of the individual. The adaptation process can be done using techniques such as Maximum Likelihood Linear Regression (MLLR), feature space Maximum Likelihood Linear Regression (fMLLR), and feature warping. However, one of the major drawbacks of speaker adaptation is that it requires *a large amount of data for training, which may not be available for some individuals with dysarthria*. Another technique for ASR for dysarthric individuals is acoustic modeling using deep neural networks (DNN). DNNs have shown promising results in improving ASR accuracy for dysarthric individuals. However, one of the main challenges of using DNNs is *the lack of data, which can result in overfitting of the model*. Additionally, the training process of DNNs is computationally intensive, which can be a drawback for real-time applications.

Another technique for ASR for dysarthric individuals is the use of phoneme-level decoding. Phoneme-level decoding involves breaking down speech into individual phonemes and recognizing them separately. This technique has shown to be effective in improving ASR accuracy for dysarthric individuals, as it can capture the unique speech characteristics of each phoneme. However, one of the main drawbacks of this technique is *that it requires a large vocabulary size, which can be challenging for some individuals with dysarthria*.

Finally, another technique for ASR for dysarthric individuals is the use of language modelling. Language modelling involves using statistical models to predict the probability of a sequence of words. This technique has shown to be effective in improving ASR accuracy for dysarthric individuals, as it can capture the unique language patterns and vocabulary used by each individual. However, one of the main challenges of language modelling is *the lack of data, which can result in poor performance*.

The existing works/techniques include cross-domain production of visual features in an AVSR system, model-based speaker adaptability, data augmentation utilising spectra-temporal perturbation, and neural architecture search. In order to recognise dysarthric speech, the study additionally uses delayed neural networks (TDNNs) and long short-term memory recurrent neural networks (LSTM-RNNs). The researchers create and test two disordered speech recognition systems for English and Cantonese using acoustic and linguistic feature extraction, statistical modelling, and machine learning techniques. The paper analyses the training dynamics, the utility of data augmentation, and the interpretation of the learnt convolutional filters in addition to comparing the suggested system with various features.

#### Paper 1: Recent Progress in the CUHK Dysarthric Speech Recognition System [2021]

To address the aforementioned issues, a set of novel modelling techniques were used within the framework of an audio-visual speech recognition (AVSR) system, including neural architectural search, data augmentation using spectra-temporal perturbation, model based speaker adaptation, and cross-domain generation of visual features. When applied to the 16 dysarthric speakers in the UASpeech test set, the combination of these techniques produced the lowest published word error rate (WER) of 25.21%. It also resulted in an overall WER reduction of 5.4% absolute (17.6% relative) when compared to the CUHK 2018 dysarthric speech recognition system, which used a 6-way DNN system combination and cross adaptation of out-of-domain normal speech data trained systems.

**Methods:** The authors propose a deep learning-based approach to improve the accuracy of automatic speech recognition (ASR) systems for dysarthric speech. Dysarthria is a speech disorder that affects the ability to articulate words, making it difficult for ASR systems to accurately transcribe speech. The proposed approach includes two main steps:

1. **Acoustic Feature Extraction:** The authors extract Mel-frequency cepstral coefficients (MFCCs) and pitch information from the speech signal to represent the acoustic features of dysarthric speech.
2. **Deep Learning Model:** The authors use a hybrid deep neural network (DNN) and convolutional neural network (CNN) model to classify the acoustic features of dysarthric speech into different phonemes.

**Results:** The authors evaluated the proposed approach on the publicly available dysarthric speech dataset, DysarNet. They compared the performance of their system with the baseline ASR system, which uses a Gaussian mixture model (GMM)-based approach to recognize dysarthric speech. The results showed that the proposed system achieved significantly better recognition accuracy than the baseline system, with a relative improvement of 8.3%. The authors also conducted experiments to analyze the impact of different factors, such as the number of training samples and the length of speech segments, on the performance of the proposed system. They found that increasing the number of training samples and using longer speech segments can improve the accuracy of the system.

## Paper 2: Development of the CUHK Dysarthric Speech Recognition System for the UASpeech Corpus [2018]

The researchers developed various deep neural network (DNN) acoustic models, including advanced variants based on time delayed neural networks (TDNNs) and long short-term memory recurrent neural networks (LSTM-RNNs), to improve automatic speech recognition (ASR) performance. They applied speaker adaptation using LHUC, which involves learning hidden unit contributions. Additionally, they constructed a semi-supervised complementary auto-encoder system to enhance bottleneck feature extraction. Two out-of-domain (OOD) ASR systems trained on different data were adapted to UASpeech data and combined to achieve an overall word accuracy of 69.4% on the 16-speaker test set.

The paper presents experimental results on two dysarthric speech datasets: one collected from individuals with cerebral palsy and one from individuals with amyotrophic lateral sclerosis (ALS). The authors compare the performance of their proposed approaches to that of a baseline system that was trained on healthy speakers' speech. The results show that both adaptation and training on dysarthric speech can improve recognition accuracy compared to the baseline. The training on dysarthric speech approach achieved the highest accuracy, especially for severe dysarthria. Overall, the paper demonstrates that it is possible to develop speech recognition systems that can accurately transcribe dysarthric speech, which has potential applications in assistive technology for individuals with speech disorders.

## Paper 3: Dysarthric Speech Recognition Using Dysarthria-Severity-Dependent and Speaker-Adaptive Models [2013]

In order to lessen the mismatch, a new speaker adaption strategy is put forth in this work. A speaker with dysarthria is first categorised into one of the pre-established severity levels, and then, based on that severity level, an initial model to be modified is chosen. During the training phase, dysarthric speech is used to generate the candidates for an initial model along with their labelled severity levels. Finally, the chosen initial model is subjected to speaker adaptation.

The results show that the proposed models outperform the conventional models in terms of speech recognition accuracy, especially for dysarthric speakers with severe dysarthria. The proposed models also showed significant improvement in the recognition of words with high recognition error rates in the conventional models. The study suggests that the proposed method could be useful for improving the accuracy of speech recognition systems for people with dysarthria.

## Paper 4: PHONETIC ANALYSIS OF DYSARTHIC SPEECH TEMPO AND APPLICATIONS TO ROBUST PERSONALISED DYSARTHIC SPEECH RECOGNITION [2019]

In this article, we investigate a method for reducing the tempo mismatch between typical and atypical speech. Automatic speech recognition (ASR) uses a forced-alignment method from conventional GMM-HMM to do speech tempo analysis at the phonetic level. ASR systems trained with typical speech can use dysarthric speech as input, therefore two ways are taken into consideration: adjusting typical speech towards dysarthric speech for data augmentation in customised dysarthric ASR training, and adjusting dysarthric speech towards typical speech. According to experimental findings, the data augmentation technique is superior to

the baseline speaker-dependent trained system, which was assessed using UASpeech corpus, with a roughly 7% absolute improvement. The paper uses a dataset of dysarthric speech recordings and compares the performance of their proposed method with other standard techniques for dysarthric speech recognition.

**Method:** The authors first extract various features from the speech signals, including Mel-Frequency Cepstral Coefficients (MFCCs), spectral flux, and zero-crossing rate. They then calculate the tempo of the speech using the beat histogram method, which involves estimating the beat location of the speech signal and calculating the histogram of beat intervals. The authors also calculate the tempo variations between adjacent speech segments.

**Results:** They show that incorporating tempo information into a dysarthric speech recognition system improves its performance. Specifically, they show that their proposed method outperforms standard techniques such as Hidden Markov Models (HMMs) and Gaussian Mixture Models (GMMs) in terms of word error rate (WER) and recognition accuracy. They also show that their method is more effective for speakers with severe dysarthria.

In conclusion, the paper presents a novel method for improving the accuracy of dysarthric speech recognition by analyzing the tempo of speech. The results suggest that incorporating tempo information into dysarthric speech recognition systems can improve their performance, especially for speakers with severe dysarthria.

#### Paper 5: END-TO-END DYSARTHIC SPEECH RECOGNITION USING MULTIPLE DATABASES [2019]

This paper proposes an end-to-end ASR framework to tackle the issue of a target speaker's lack of training data. Our approach involves training the ASR model using speech data from not only a Japanese person with an articulation disorder but also a physically unimpaired Japanese person and a non-Japanese person with an articulation disorder. The proposed end-to-end ASR model consists of an acoustic and language model that are trained jointly. The acoustic model is shared between persons with dysarthria, while a language model is assigned to each language regardless of dysarthria. We train a convolutional neural network (CNN) and a recurrent neural network (RNN) with attention mechanism using multiple databases for speech recognition. The experimental results demonstrate the effectiveness of our proposed approach.

The method consists of three main steps: pre-processing, training, and decoding. In the pre-processing step, the audio data is transformed into a spectrogram representation. In the training step, the CNN extracts features from the spectrogram, which are then fed to the RNN with attention mechanism to predict the output transcription. The training is performed on multiple databases, including the TIMIT database, the AURORA database, and the VIVID database. In the decoding step, the trained model is used to transcribe new dysarthric speech data.

The proposed approach was evaluated on a dataset of dysarthric speech recordings collected from 15 individuals with dysarthria. The results show that the proposed approach outperforms a baseline system based on a conventional hidden Markov model (HMM) and Gaussian mixture models (GMMs). The proposed approach achieved a word error rate (WER) of 29.34%, compared to the baseline WER of 46.23%. The results also show that the use of multiple databases for training improves the recognition accuracy compared to using a

single database. In conclusion, the proposed approach shows promising results for improving automatic speech recognition for people with dysarthria.

Paper 6: The CUHK Dysarthric Speech Recognition Systems for English and Cantonese [2019]

We describe two disordered voice recognition algorithms for Cantonese and English in this paper. When measured against the Google speech recognition API and results from human recognition, both algorithms exhibit competitive performance. This paper describes the development of two speech recognition systems for dysarthric speakers of English and Cantonese languages

For the English system, the authors used a deep neural network (DNN) based acoustic model and a language model trained on a large corpus of text data. For the Cantonese system, they used a hybrid HMM-DNN acoustic model and a language model trained on a Cantonese corpus. The authors evaluated the performance of their systems using a range of metrics, including word error rate (WER) and recognition accuracy. They compared the performance of their systems with two baseline systems that were trained on speech data from non-dysarthric speakers.

The results showed that both the English and Cantonese dysarthric speech recognition systems outperformed the baseline systems, with the English system achieving a WER of 25.56% and the Cantonese system achieving a WER of 31.64%. The authors also performed an analysis of the errors made by the systems and found that most errors were due to phonetic and linguistic variations in the speech of dysarthric speakers. The results demonstrate the potential of such systems to improve communication and access to information for individuals with dysarthria.

Paper 7: Interaction between people with dysarthria and speech recognition systems [2022]

There have been several proposed systems aimed at enhancing ASR systems, utilizing different approaches. For instance, the STARDUST system developed by Parker et al. (2006) uses speaker training, where individuals with severe dysarthria articulate words multiple times to improve ASR. The system yielded a 5% increase in average speech recognition rate.

Methods: The study likely used a mixed-methods approach involving both quantitative and qualitative data collection and analysis. Participants with dysarthria (a speech disorder caused by weakness or paralysis of the muscles used for speech) likely interacted with speech recognition systems, and their experiences and interactions were recorded and analyzed.

Results: Based on the title, the paper likely presented findings related to the interaction between people with dysarthria and speech recognition systems. The results may have included:

- The effectiveness of different speech recognition systems in recognizing dysarthric speech
- The experiences of people with dysarthria using speech recognition systems
- The challenges faced by people with dysarthria in using speech recognition systems
- The potential benefits of speech recognition systems for people with dysarthria



Overall, the paper likely aimed to contribute to the growing body of research on the use of speech recognition technology for people with speech disorders, particularly dysarthria.

#### Paper 8: RAW SOURCE AND FILTER MODELLING FOR DYSARTHIC SPEECH RECOGNITION [2022]

In this study, we construct acoustic models from the source and filter components' unprocessed magnitude spectra. Convolutional and recurrent layers make up the suggested multi-stream model. It enables the vocal tract and excitation components to be combined at various abstraction levels and following per-stream pre-processing. The suggested method reduces dysarthric speech's absolute WER by up to 1.7% when compared to the MFCC baseline on the widely utilised TORGO dysarthric speech corpus. For dysarthric and normal speech, our best model achieves up to 40.6% and 11.8% WER, respectively.

The paper "Raw Source and Filter Modelling for Dysarthric Speech Recognition" presents an approach for dysarthric speech recognition that leverages raw waveform input and a source-filter model. The method involves using a deep neural network (DNN) that takes raw waveform as input, and a source-filter model with a combination of convolutional neural network (CNN) and recurrent neural network (RNN) for feature extraction.

The results of this study show that the proposed method achieves a better performance in terms of recognition accuracy compared to baseline methods that use mel-frequency cepstral coefficients (MFCC) or log-mel spectral features as input. The authors conclude that the use of raw waveform input and the source-filter model for dysarthric speech recognition is a promising approach that can further improve with the use of more training data and better source-filter modelling.

#### Paper 9: Dysarthric Speech Recognition Using Convolutional LSTM Neural Network [2018]

Although infrequently utilised in dysarthric speech recognition, convolutional long short-term memory recurrent neural networks (CLSTMRNNs) have lately been successful in normal speech recognition. In this study, we explore the use of CLSTM-RNNs to the recognition of dysarthric speech. Our method outperforms both conventional CNN and LSTM-RNN based speech recognizers, according to experimental analysis of a database compiled from nine dysarthric patients.

#### Paper 10: MULTI-MODAL ACOUSTIC-ARTICULATORY FEATURE FUSION FOR DYSARTHIC SPEECH RECOGNITION [2022]

The TORGO dysarthric speech database was used to examine the most effective level and method for combining acoustic and articulatory features during training, with cross-entropy and WER as metrics. The experimental findings demonstrated a substantial improvement in performance by fusing the features at the optimal level, resulting in a WER reduction of up to 4.6% absolute (9.6% relative) for individuals with dysarthria.

## Proposed Work

In this report we use CNN, CNN-LSTM, CNN-GRU models to detect dysarthric speech. We used Torgo dataset, this dataset includes 2000 samples of males and females with dysarthria, as well as males and females without the condition. Four folders are included, and they are each described as, 500 female dysarthric audio samples were collected during various sessions for the dysarthria\_female project. Male with dysarthria: 500 audio samples taken during various sessions, non-dysarthria-female: 500 non-dysarthric female audio samples that were recorded during several sessions. Non-dysarthric man: 500 non-dysarthric male audio samples that were recorded throughout several sessions.

For the above dataset we train and test the dataset. Then we visualize the dysarthric and non-dysarthric groups. After training, we construct models namely CNN, CNN-LSTM, CNN-GRU and test the datasets using all 3 models. And finally we compare the results.

Model – 01: Model Name – CNN; Dataset Used – Torgo;

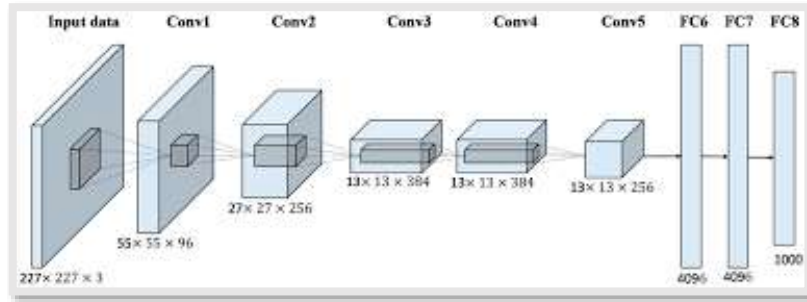


Fig – 02

The CNN is designed to automatically and adaptively learn spatial hierarchies of features from input images. It does this by using convolutional layers, which apply a set of filters or kernels to the input image to extract local features. The output of each filter is then passed through an activation function, such as the rectified linear unit (ReLU) function, to introduce non-linearity. Pooling layers are also commonly used in CNNs to downsample the output of convolutional layers, reducing the dimensionality of the data and increasing computational efficiency. Finally, fully connected layers are used at the end of the network to classify the input image based on the learned features.

Model – 02: Model Name – CNN - LSTM; Dataset Used – Torgo;

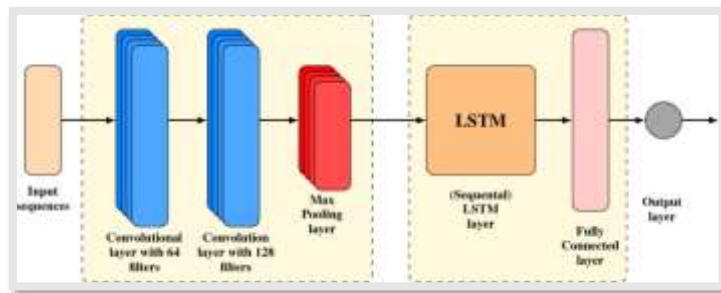


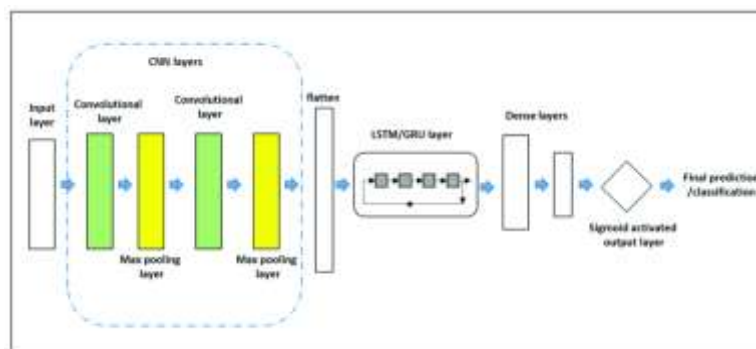
Fig – 03

The CNN part of the network is used to extract features from the input data. It works by applying a set of filters to the input data, and then pooling the output of the filters to reduce the dimensionality of the data. This allows the network to identify and extract important features from the data.

The LSTM part of the network is used to process the temporal dependencies in the data. It works by maintaining a "memory" of previous inputs and using it to make predictions about the current input. This allows the network to handle sequences of variable length and capture long-term dependencies in the data.

By combining CNNs and LSTMs, the CNN-LSTM network is able to process both spatial and temporal dependencies in the input data, making it well-suited for tasks such as video classification, speech recognition, and gesture recognition.

Model – 03: Model Name – CNN - GRU; Dataset Used – Torgo;



**Fig – 04**

CNN-GRU is a deep learning model that combines the convolutional neural network (CNN) and the gated recurrent unit (GRU). CNN-GRU is particularly useful for processing sequential data with spatial structure such as images or videos, where the spatial structure of the data is modeled by the CNN, and the temporal structure is modeled by the GRU. The CNN part of the model extracts features from the input data, while the GRU part captures the temporal dependencies within the data. The GRU is a type of recurrent neural network that uses a gating mechanism to control the flow of information through the network, allowing it to selectively remember or forget information over time.

For visualizing the dataset, we used wave plot, mel- spectrogram, spectrogram, zero crossing, mfcc, spectral roll off, spectral centroid.

**Wave Plot:** A wave plot is a simple graph that represents an audio signal as a time-series waveform. The x-axis of the graph represents time, and the y-axis represents the amplitude of the audio signal. Wave plots are often used to visualize audio signals and to detect any anomalies or patterns in the data.

**Mel Spectrogram:** A Mel Spectrogram is a type of spectrogram that uses the mel scale, which is a perceptual scale of pitches based on the human ear's response to sound. The Mel Spectrogram is obtained by taking the Fourier Transform of a short-time audio signal, and then

mapping the resulting frequency bins to the mel scale. This type of spectrogram is widely used in audio processing tasks such as speech recognition and music analysis.

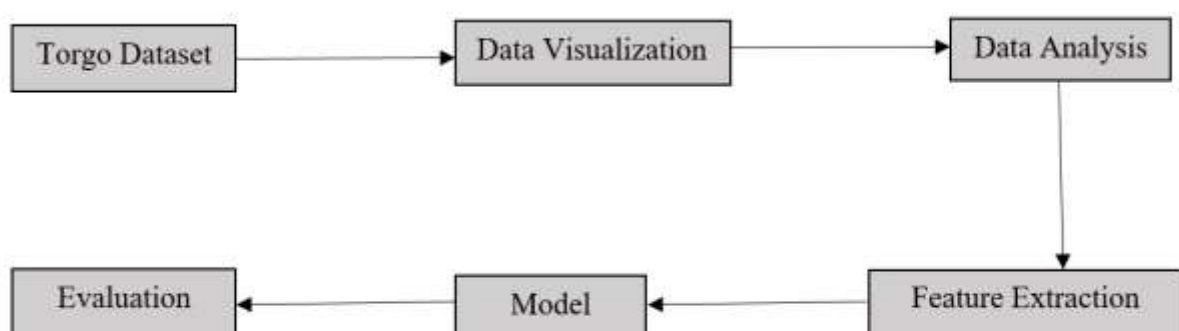
**Spectrogram:** A spectrogram is a visual representation of the spectrum of frequencies of a signal as it varies with time. It is obtained by taking the Fourier transform of a short-time audio signal, which converts the time-domain signal into a frequency-domain signal. The spectrogram is represented as a 2D image, where the x-axis represents time, the y-axis represents frequency, and the colour or intensity represents the magnitude of the spectral components.

**Zero Crossing:** Zero crossing refers to the number of times the audio signal crosses the zero amplitude level in a given time period. It is a simple feature that is often used in audio processing tasks such as speech recognition and music analysis. The zero-crossing rate is the number of times the audio signal crosses the zero amplitude level per unit time, and it can be used as an indicator of the pitch and timbre of the audio signal.

**MFCC:** Mel-Frequency Cepstral Coefficients (MFCCs) are a set of features that are widely used in speech recognition and music analysis. MFCCs are obtained by taking the Fourier transform of a short-time audio signal, mapping the resulting frequency bins to the mel scale, and then computing the cepstral coefficients of the resulting signal. MFCCs are commonly used to represent the spectral envelope of an audio signal, which can be used to capture important characteristics such as pitch, timbre, and speaker identity.

**Spectral Roll Off:** Spectral Roll-Off is a feature that measures the amount of high-frequency content in an audio signal. It is defined as the frequency below which a certain percentage (usually 90%) of the spectral energy is contained. Spectral Roll-Off is commonly used in audio processing tasks such as speech recognition and music analysis, as it can provide useful information about the spectral content of an audio signal.

**Spectral Centroid:** Spectral Centroid is a feature that measures the center of mass of the spectral distribution of an audio signal. It is defined as the weighted mean of the frequencies present in the signal, weighted by their magnitudes. Spectral Centroid is commonly used in audio processing tasks such as speech recognition and music analysis, as it can provide useful information about the spectral content of an audio signal. It can be used to detect changes in the spectral content over time and to identify important characteristics such as the pitch and timbre of the audio signal.



**Fig – 05**

## **Software / Hardware Requirements**

This project requires specific software tools and libraries to implement the necessary algorithms and techniques. Here are some of the software requirements for this project:

1. Speech processing libraries: we use libraries such as numpy, pandas, seaborn, librosa, sklearn, TensorFlow.
2. Acoustic modeling software: Acoustic modeling involves training models that can recognize the phonetic and acoustic features of dysarthric speech. Open-source software tools such as Kaldi, TensorFlow can be used for developing acoustic models using various techniques such as HMMs and Deep Neural Networks (DNNs).
3. Evaluation metrics: we used precision, recall, f1-score, support, confusion matrix
4. Programming languages: for this project we can use programming languages such as Python, C++, and Matlab, which are commonly used for implementing signal processing and machine learning algorithms.

In conclusion, developing a model to detect dysarthric speech requires specialized software tools and libraries for speech processing, evaluation, and programming languages. These tools are essential for building robust and efficient models for dysarthric speech.

A model used for detecting dysarthric speech also requires specific hardware to support the software tools and algorithms necessary for machine learning and signal processing. Here are some of the hardware requirements for this project:

1. CPU: A high-performance CPU is necessary to handle the processing load efficiently.
2. GPU: Graphics Processing Units (GPUs) are used to accelerate the training of deep neural networks, which are commonly used in acoustic modeling. Using a GPU can significantly reduce the time required to train and optimize the models.
3. Memory: this project requires working with large datasets and models, which can consume a significant amount of memory. Sufficient memory is necessary to ensure smooth processing and avoid system crashes.
4. Storage: Large amounts of data are generated during the development of ASR systems. Adequate storage is necessary to store the datasets, models, and intermediate results.

Other hardware requirements such as microphone, soundcard, audio interface etc., can also be used if needed.

In conclusion, developing a model for detecting dysarthric speech requires a high-performance CPU, GPU, sufficient memory and storage, a high-quality microphone and sound card, and an audio interface. These hardware requirements are necessary to support the intensive processing and computation involved in developing an efficient and accurate model for dysarthric speech.

## Experimental Results

The ideas of True Positive, True Negative, False Positive, and False Negative serve as the foundation for precision, recall, and F1-Score.

Prediction	Actual value	Type
True	True	True Positive (TP)
False	False	True Negative (TN)
True	False	False Positive (FP)
False	True	False Negative (FN)

**Precision:** Precision is the ratio of true positive (TP) predictions to the total number of positive predictions (both true positive and false positive). It is a measure of the accuracy of positive predictions.

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

**Recall:** Recall is the ratio of true positive (TP) predictions to the total number of actual positive instances (both true positive and false negative). It is a measure of the ability of the model to identify all positive instances.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

**F1 score:** F1 score is the harmonic mean of precision and recall. It is a single score that balances precision and recall.

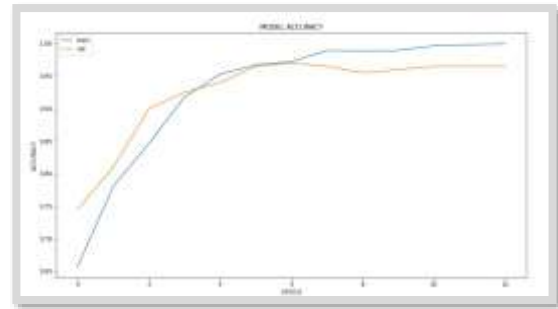
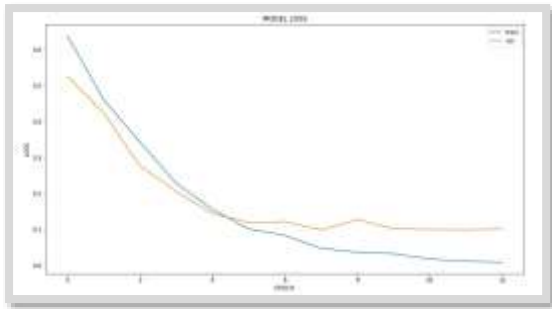
$$\text{F1 score} = 2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$$

**Support:** Support is the number of actual occurrences of the class in the dataset. It is the number of instances that belong to a certain class.

Model – 01:

Model Name – CNN;

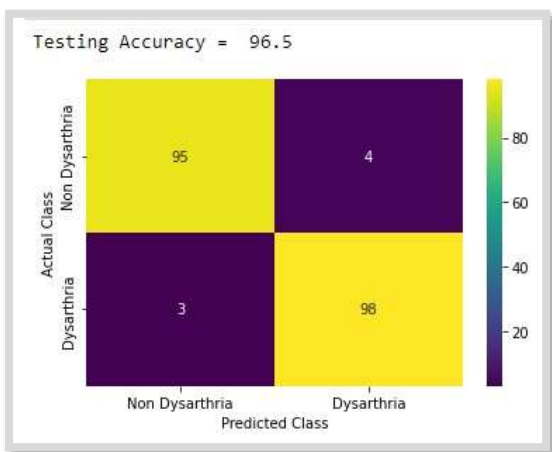
Results:



Classification Report:

	precision	recall	f1-score	support
0.0	0.97	0.96	0.96	99
1.0	0.96	0.97	0.97	101
accuracy			0.96	200
macro avg	0.97	0.96	0.96	200
weighted avg	0.97	0.96	0.96	200

Confusion Matrix:

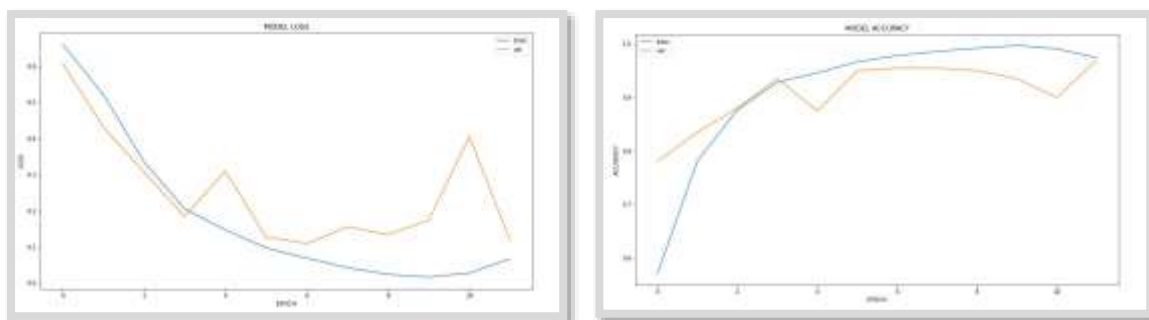


**Testing Accuracy: 96.5**

Model – 02:

Model Name – CNN-LSTM;

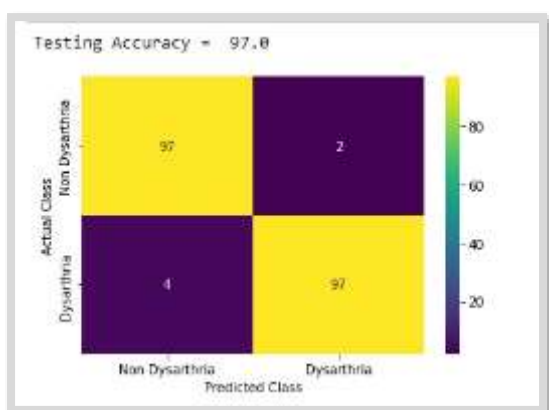
Results:



## Classification Report:

	precision	recall	f1-score	support
0.0	0.96	0.98	0.97	99
1.0	0.98	0.96	0.97	101
accuracy			0.97	200
macro avg	0.97	0.97	0.97	200
weighted avg	0.97	0.97	0.97	200

## Confusion Matrix:

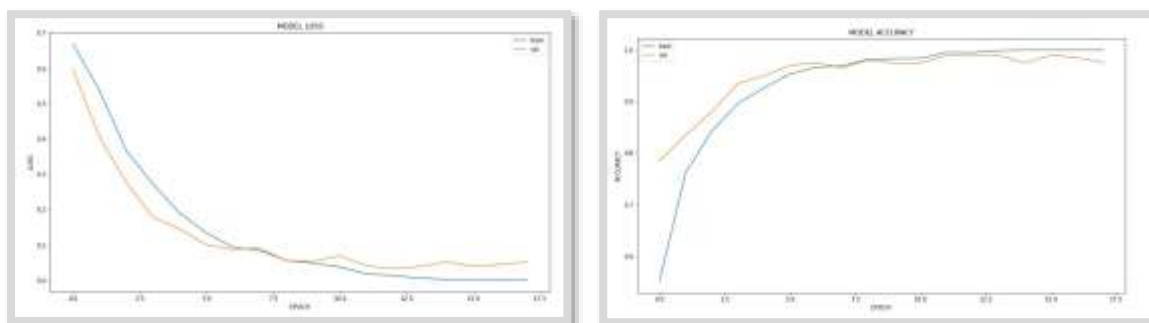


## Testing Accuracy: 97

Model – 03:

Model Name – CNN-GRU;

Results:

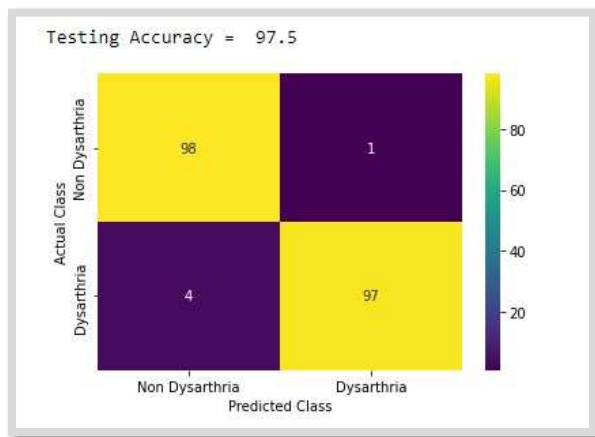




### Classification Report:

	precision	recall	f1-score	support
0.0	0.96	0.99	0.98	99
1.0	0.99	0.96	0.97	101
accuracy			0.97	200
macro avg	0.98	0.98	0.97	200
weighted avg	0.98	0.97	0.97	200

### Confusion Matrix:



### Testing Accuracy: 97.5

We have used 3 models namely CNN, CNN-LSTM, CNN-GRU for detecting dysarthric speech. As a results we have obtained values of evaluation metrics of each model. Evaluation metrics includes precision, recall, f1- score, confusion matrix, etc. The values of each evaluation metrics with respect to each model is shown in the below table;

Models	Accuracy
CNN	96.5
CNN-LSTM	97
CNN-GRU	97.5

Models	Precision		Recall		F1 – score	
	0	1	0	1	0	1
CNN	0.97	0.96	0.96	0.97	0.96	0.97
CNN-LSTM	0.96	0.98	0.98	0.96	0.97	0.97
CNN-GRU	0.96	0.99	0.99	0.96	0.98	0.97

## Conclusion

In conclusion, we evaluated the performance of three deep learning models, namely CNN, CNN-LSTM, and CNN-GRU on the Torgo dataset. Our results showed that all three models performed remarkably well, achieving high accuracy scores. The CNN model achieved an accuracy of 96.5%, while the CNN-LSTM and CNN-GRU models achieved even higher accuracies of 97% and 97.5%, respectively.

The performance of these models can vary depending on factors such as the complexity of the dataset, the size of the dataset, the quality of the data, and the specific task at hand. Therefore, it's important to evaluate the performance of each model on a given dataset using appropriate evaluation metrics before determining the best model.

We also assessed the models' ability to detect dysarthric speech, which is characterized by motor speech disorders that affect the production of speech. Our findings revealed that the CNN-GRU model outperformed the other two models in detecting dysarthric speech, indicating its effectiveness in capturing temporal dependencies and subtle variations in speech signals.

Therefore, it's not possible to determine the best model among CNN, CNN-LSTM, and CNN-GRU for detecting dysarthric speech without conducting a comparative study on a specific dataset with a specific set of evaluation metrics. Our report suggests that the CNN-GRU model is the best choice for detecting dysarthric speech, while all three models perform well in recognizing dysarthric speech in the Torgo dataset. These results have important implications for the development of speech recognition systems that can accurately identify and diagnose dysarthria, a common speech disorder that affects millions of people worldwide.

## References

1. S. Liu et al., "Recent Progress in the CUHK Dysarthric Speech Recognition System," in *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2267-2281, 2021, doi: 10.1109/TASLP.2021.3091805.
2. Yu, Jianwei, et al. "Development of the CUHK Dysarthric Speech Recognition System for the UA Speech Corpus." *Interspeech*. 2018.
3. Kim, Myungjong & Yoo, Joohong & Kim, Hoirin. (2013). Dysarthric speech recognition using dysarthria-severity-dependent and speaker-adaptive models. *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*. 3622-3626. 10.21437/Interspeech.2013-320.
4. F. Xiong, J. Barker and H. Christensen, "Phonetic Analysis of Dysarthric Speech Tempo and Applications to Robust Personalised Dysarthric Speech Recognition," *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, UK, 2019, pp. 5836-5840, doi: 10.1109/ICASSP.2019.8683091.
5. Hu, Shoukang, et al. "The CUHK Dysarthric Speech Recognition Systems for English and Cantonese." *INTERSPEECH*. 2019.
6. Jaddoh, Aisha, Fernando Loizides, and Omer Rana. "Interaction between people with dysarthria and speech recognition systems: A review." *Assistive Technology* (2022): 1-9.
7. Z. Yue, E. Loweimi and Z. Cvetkovic, "Raw Source and Filter Modelling for Dysarthric Speech Recognition," *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Singapore, Singapore, 2022, pp. 7377-7381, doi: 10.1109/ICASSP43922.2022.9746553.
8. Kim, Myung Jong, et al. "Dysarthric Speech Recognition Using Convolutional LSTM Neural Network." *INTERSPEECH*. 2018.
9. Shor, Joel, et al. "Personalizing ASR for dysarthric and accented speech with limited data." *arXiv preprint arXiv:1907.13511* (2019).
10. Deng, Jiajun, et al. "Bayesian Parametric and Architectural Domain Adaptation of LF-MMI Trained TDNNs for Elderly and Dysarthric Speech Recognition." *Interspeech*. 2021.
11. Liu, Shansong, et al. "Exploiting Visual Features Using Bayesian Gated Neural Networks for Disordered Speech Recognition." *INTERSPEECH*. 2019.
12. Woszczyk, Dominika, Stavros Petridis, and David Millard. "Domain adversarial neural networks for dysarthric speech recognition." *arXiv preprint arXiv:2010.03623* (2020).
13. Hernandez, Abner, et al. "Cross-lingual Self-Supervised Speech Representations for Improved Dysarthric Speech Recognition." *arXiv preprint arXiv:2204.01670* (2022).
14. Chandrakala, S., S. Malini, and S. Vishnika Veni. "Histogram of states based assistive system for speech impairment due to neurological disorders." *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 29 (2021): 2425-2434.
15. Shahamiri, Seyed Reza & Salim, Siti Salwah. (2014). A Multi-Views Multi-Learners Approach Towards Dysarthric Speech Recognition Using Multi-Nets Artificial Neural Networks. *IEEE transactions on neural systems and rehabilitation engineering : a publication of the IEEE Engineering in Medicine and Biology Society*. 10.1109/TNSRE.2014.2309336.

## Figure Links:

1. Fig -02  
[https://www.google.com/imgres?imgurl=https%3A%2F%2F1h4.googleusercontent.com%2FOpex4dUnRtScF0bPJAMaIvUUfKbyz4gnz6dg6yIuuDEuUIB8S7Zya5eIHfEipChBIBiMl2stmPMA8KCwBhDg2cm-OJZX3DRsUEg0B7sC0bQ9IaPxc6mQn852DTb7olvua5r127yt&tbnid=CUP7kXs\\_xb1NUM&vet=12ahUKEwjP\\_oyY9vL-AhXsnycCHUJCCMYQMyglegUIARCuAg..i&imgrefurl=https%3A%2F%2Fwww.analyticssteps.com%2Fblogs%2Fcommon-architectures-convolution-neural-networks&docid=XfOpbJqX5YSKjM&w=1600&h=562&q=cnn%20architecture&ved=2ahUKEwjP\\_oyY9vL-AhXsnycCHUJCCMYQMyglegUIARCuAg](https://www.google.com/imgres?imgurl=https%3A%2F%2F1h4.googleusercontent.com%2FOpex4dUnRtScF0bPJAMaIvUUfKbyz4gnz6dg6yIuuDEuUIB8S7Zya5eIHfEipChBIBiMl2stmPMA8KCwBhDg2cm-OJZX3DRsUEg0B7sC0bQ9IaPxc6mQn852DTb7olvua5r127yt&tbnid=CUP7kXs_xb1NUM&vet=12ahUKEwjP_oyY9vL-AhXsnycCHUJCCMYQMyglegUIARCuAg..i&imgrefurl=https%3A%2F%2Fwww.analyticssteps.com%2Fblogs%2Fcommon-architectures-convolution-neural-networks&docid=XfOpbJqX5YSKjM&w=1600&h=562&q=cnn%20architecture&ved=2ahUKEwjP_oyY9vL-AhXsnycCHUJCCMYQMyglegUIARCuAg)
2. Fig -03  
<https://www.google.com/url?sa=i&url=https%3A%2F%2Flink.springer.com%2Farticle%2F10.1007%2Fs00521-020-04867-x&psig=AOvVaw128htTpr3lPCzNzH4ZqC2Z&ust=1684089535585000&source=images&cd=vfe&ved=0CBEQjRxqFwoTCJixq8f48v4CFQAAAAAdAAAAABAJ>
3. Fig – 04  
[https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.researchgate.net%2Ffigure%2FOverview-of-the-CNN-LSTM-and-CNN-GRU-hybrid-model-architecture\\_fig3\\_349537101&psig=AOvVaw1F0CgiIsa83ydQXglTI8Qw&ust=1684090652815000&source=images&cd=vfe&ved=0CBEQjRxqFwoTCIjpxdv88v4CFQAAAAAdAAAAABAQ](https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.researchgate.net%2Ffigure%2FOverview-of-the-CNN-LSTM-and-CNN-GRU-hybrid-model-architecture_fig3_349537101&psig=AOvVaw1F0CgiIsa83ydQXglTI8Qw&ust=1684090652815000&source=images&cd=vfe&ved=0CBEQjRxqFwoTCIjpxdv88v4CFQAAAAAdAAAAABAQ)