# Classification of Dysarthria based on the Levels of Severity. A Systematic Review

Afnan Al-Ali<sup>a</sup>, Somaya Al-Maadeed<sup>a</sup>, Moutaz Saleh<sup>a</sup>, Rani Chinnappa Naidu<sup>b</sup>, Zachariah C Alex<sup>b</sup>, Prakash Ramachandran<sup>b</sup>, Rajeev Khoodeeram<sup>c</sup> and Rajesh Kumar M<sup>b</sup>

## ARTICLE INFO

Keywords:
Dysarthria,
Classification,
Severity Levels,
Artificial Intelligence (AI)-based models,
Intelligibility

#### ABSTRACT

Dysarthria is a neurological speech disorder that can significantly impact affected individuals' communication abilities and overall quality of life. The accurate and objective classification of dysarthria and the determination of its severity are crucial for effective therapeutic intervention. While traditional assessments by speech-language pathologists (SLPs) are common, they are often subjective, time-consuming, and can vary between practitioners. Emerging machine learning-based models have shown the potential to provide a more objective dysarthria assessment, enhancing diagnostic accuracy and reliability. This systematic review aims to comprehensively analyze current methodologies for classifying dysarthria based on severity levels. Specifically, this review will focus on determining the most effective set and type of features that can be used for automatic patient classification and evaluating the best AI techniques for this purpose. We will systematically review the literature on the automatic classification of dysarthria severity levels. Sources of information will include electronic databases and grey literature. Selection criteria will be established based on relevance to the research questions. Data extraction will include methodologies used, the type of features extracted for classification, and AI techniques employed. The findings of this systematic review will contribute to the current understanding of dysarthria classification, inform future research, and support the development of improved diagnostic tools. The implications of these findings could be significant in advancing patient care and improving therapeutic outcomes for individuals affected by dysarthria.

# 1. Introduction

Speech is a distinctive, intricate, dynamic motor activity that enables us to articulate our thoughts and emotions and interact with and regulate our surroundings. Furthermore, it constitutes one of the five heritable verbal traits, encompassing speech, language, reading, writing, and spelling. Speech involves integrating neurocognitive processes for organizing thoughts into language, motor speech planning for executing verbal messages, and neuromuscular execution for coordinating speech muscles. Together, these processes constitute motor speech activities.

When neurologic impairments affect these motor speech activities, a speech disorder will result, which can also be known as Motor speech disorder (MSD) [1][2].

There are two main types of MSD: dysarthrias and apraxia of speech. Dysarthria is a group of neurologic speech disorders involving abnormalities in speech production's movement aspects. These disorders manifest as changes in strength, speed, range, steadiness, tone, or accuracy of movements required for breathing, phonation, resonance, articulation, or prosody. Sensorimotor abnormalities, such as weakness, spasticity, incoordination, involuntary movements, or variations in muscle tone, underlie dysarthria. The condition is specifically neurologic and can be categorized into distinct types based on perceptual characteristics and underlying neuropathophysiology. It is essential to define dysarthria accurately to differentiate it from other speech and language disorders

aa1805360@qu.edu.qa (A. Al-Ali)
ORCID(s):

and to ensure its meaningful application in research and clinical settings [2]. While Verbal apraxia, which is also known as apraxia of speech, is a contentious condition that some view as a deficit in the motor planning of speech. This disorder is marked by "higher order" errors, including metathesis and segment addition, alongside errors that suggest a lack of coordination in articulation. These characteristics suggest a relatively significant level of damage to the neural system [3].

Traditionally, dysarthria is classified into several subtypes, including spastic, flaccid, hypokinetic, hyperkinetic, ataxic, and mixed. More recently, additional subtypes such as unilateral upper motor neuron dysarthria and undetermined dysarthria have been recognized [4]. Dysarthria can occur at any stage of life. Common causes encompass stroke, severe head injury, brain tumors, Parkinson's disease, multiple sclerosis, motor neuron disease, cerebral palsy, Down's syndrome, and certain medications, including those used to treat epilepsy, which may induce a side effect [2]. The severity of dysarthria is determined by the degree of involvement in the affected body regions caused by the underlying condition. Assessment of dysarthria means mainly grading its severity. This is typically performed by speechlanguage pathologists (SLPs) using some descriptive terms, but it can be a time-consuming and labor-intensive process with variations between different SLPs. Thus, it is essential to have objective methods of evaluating the level of intelligibility in dysarthria cases [5]. For this purpose, machine learning-based models were developed for automatic assessment of dysarthrias' levels of severity to achieve enhanced diagnostic accuracy, consistency, and reliability, all while maintaining cost-effectiveness and

<sup>&</sup>lt;sup>a</sup> Computer Science and Engineering Department, Qatar University, Doha, Qatar,

<sup>&</sup>lt;sup>b</sup> Vellore Institute of Technology, Vellore, India,

<sup>&</sup>lt;sup>c</sup> Université Des Mascareignes, Mauritius,

<sup>\*</sup>Corresponding author

expediency [6]. In both paths, features are extracted from the candidate samples to help the underlying system for classification, where there are specific sets of features for each type of dysarthria. To our knowledge, this is the first comprehensive review that analyzes the works of classifying dysarthria cases based on the levels of severity. In the literature, several gaps warrant further research in conducting a systematic review on classifying dysarthria based on severity levels. These gaps include the absence of a review focused explicitly on severitybased classification, limited research on severity levels within specific populations [7][8][9][10] (populationspecific considerations), a lack of standardized measures for assessing severity [8][11][12][13] (measurement challenges), the integration of technology in severity classification [14] (technological advancements), a need for a comparative analysis of existing classification systems [15] (evaluation of existing approaches), specific causes of dysarthria [16] (etiological factors), general treatment helping clinicians for stable dysarthria [17] (treatment approaches) and automated intelligibility assessment [18] (technology-driven assessment methods). Addressing these gaps through a systematic review would offer valuable insights for clinicians, researchers, and stakeholders involved in dysarthria assessment and treatment

In this systematic review, we aim to answer the below research questions:

- 1. What are the optimal set and type of features that are consistently effective across various severity levels of dysarthria, enabling automated classification of patients?
- 2. Which artificial intelligence techniques are most suitable for accurately classifying dysarthria patients, considering factors such as training time efficiency and high accuracy?

The rest of the paper is organized as follows: Section 2 will show the strategy followed in our research, Section 3 presents how dysarthric cases are classified based on clinical techniques, Section 4 will show the same for Section 3, but based on machine learning techniques, Section 5 discusses the results shown in the two previous sections and analyses them. Section 6 refers to some common limitations of the related work. Section 7 suggests a few points to solve the gap in this research area, and finally, the Conclusion in Section 8 will summarize our work and highlights our findings in this review.

# 2. Search Strategy

A comprehensive exploration was conducted on various electronic databases, including ACM, EMBASE, SpringerLink, PubMed, Scopus, IEEE, MDPI, Elsevier, and some other conferences popular in the area, based on the search keywords: "classification" AND "Dysarthria" AND "severity levels" or "Assessment" AND "Dysarthria" AND "Intelligibility." Figure 1 shows the search strategy process.

A total of 978 publications were found. After deleting the duplicates, 733 publications passed through a screening process where only the articles which contained the

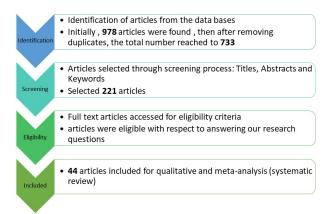


Figure 1: Search Study Process.

keywords in their titles, abstracts, and keywords were selected, and the excluded publications were based on the following criteria:(1) reviews and surveys are not considered, (2) articles do not have the main keywords in their titles and abstracts, (3) general chapters related to motor speech disorders are not considered, and (4) other types of disorders like Dysphagia or Dementia.

The second stage of the study search is knowing the eligibility of each publication to be selected for the analysis later on in this study. The exclusion criteria in this stage are (1) The tools or software for rehabilitation or therapy, (2) automatic recognition systems or the identification of dysarthria, (3) reviews and surveys related to dysarthria or one of its specific types or causes, (4) the impact of specific features of dysarthric cases, and (5) the binary classification of dysarthria.

The final publications are 44 articles mainly related to classifying dysarthria patients based on severity levels and assessing the intelligibility factor as severity levels.

As our primary goal is to classify dysarthric cases based on severity levels, we will refer to both types: clinical or human-based and AI or machine learning-based techniques. We will highlight the methods, extracted features, and the datasets they evaluated the models on. Figure 2 shows the taxonomy of our review in classifying dysarthria based on the severity levels.

# 3. Classification based on Clinical Techniques

When a dysarthric patient is admitted to the hospital, a specific procedure will be followed based on the clinical setting, available resources, and individual needs. The main points regarding each procedure are explained in the sub-sections below:

#### 3.1. Methods

As a clinical procedure, speech-language pathologists (SLPs) follow up on dysarthrias cases. Other rehabilitation professionals, physicians, and nurses may also be involved in such treatment. Each SLPs will search for some descriptive patterns in the patient to confirm his diagnosis and decide the type of treatment based on the severity of the case [19].

The procedure followed by speech-language pathologists (SLPs) to measure the severity of dysarthria in a patient typically involves employing a comprehensive approach to assess dysarthria in patients. This process includes gathering the patient's medical history to provide context, visually observing their speech and oral motor movements, making perceptual assessments of speech characteristics and severity, evaluating speech intelligibility through tests and measurements, assessing the impact of dysarthria on functional communication abilities, and collaborating with patients, family members, and healthcare professionals to gain a comprehensive understanding of the condition. By following this systematic approach, SLPs can develop tailored treatment plans to address the specific challenges faced by each patient[2] [20].

### 3.2. Features

Speech-language pathologists (SLPs) use various features to assess the severity of dysarthria in patients. These features include articulation, by assessing the accuracy and precision of speech sound production; phonation, by evaluating voice quality, pitch, loudness, and presence of abnormalities; resonance, by examining the control of the velopharyngeal mechanism for speech clarity; prosody, by analyzing rhythm, stress, intonation, and melodic contour; speech rate, by observing speed, pacing, and pauses; and intelligibility, by determining the percentage of intelligible speech.

SLPs employ clinical observation, perceptual judgments, and instrumental assessments like acoustic analysis to assess these features. By considering these aspects, SLPs can determine dysarthria severity and develop customized treatment plans to address patients' specific speech challenges [19][21].

## 3.3. Evaluation

Several standardized rating scales and perceptual judgments are commonly used to assess dysarthria's severity and specific features. These include formal and informal assessments. The formal ones are represented by the most famous Frenchay Dysarthria Assessment (FDA) [22], which evaluates respiration, phonation, articulation, and prosody; other measurements like the Assessment of Intelligibility of Dysarthric Speech (AIDS), which measures overall speech intelligibility [23] and dysarthria profile [24]; the and Voice Handicap Index [25]. Informal assessments, such as oral motor examinations [26], are often used with formal assessments. Perceptual assessment [27] is used by speech-language pathologists, relying on their expertise in active listening and analyzing speech. It's important to note that these perceptual judgments are subjective, emphasizing the need for skilled clinicians to accurately assess and interpret the speech characteristics of individuals with dysarthria. Finally, Communication Activities of Daily Living-Second Edition (CADL-2) [28], Although not specific to dysarthria, this assessment measures functional communication abilities in daily life situations. It evaluates the impact of dysarthria on communication in various contexts[29].

The inclusion of the section on clinical assessment in our review serves two essential purposes. Firstly, it provides a comprehensive understanding of the traditional approaches and methods used in dysarthria classification based on severity levels. By examining the techniques employed by clinicians and the features they consider, we gain insights into the established practices and evaluation tools that have been relied upon for decades. This knowledge is crucial for contextualizing and appreciating the advancements brought by AI-based techniques.

Secondly, incorporating the clinical assessment perspective allows for a comparative analysis between human-based approaches and AI-based methods. By juxtaposing the strengths and limitations of clinical judgment with the capabilities of AI models, we can critically evaluate the potential of AI to augment and enhance dysarthria classification. This comparative analysis enables us to identify the unique contributions of AI techniques, such as increased objectivity, scalability, and potential for automated assessment.

## 4. Classification based on AI Techniques

Despite the effectiveness of SLPs' efforts, this process is time-consuming, mainly subjective, and suffers from differences among individual opinions. This gap inspired the researchers to seek more accurate models, away from the subjective perspective, using AI-based models, as explained below.

## 4.1. Methods

Machine learning and deep learning-based models have shown promise in assessing the severity of dysarthria, offering alternative approaches to traditional assessments conducted by SLPs. These models leverage the power of artificial intelligence to analyze speech patterns and extract meaningful information for severity evaluation [30].

Several common machine learning approaches, such as support vector machines (SVM), Random Forests (RF), or Artificial Neural Networks (ANN), have been utilized to develop predictive models [31]. These models are trained on a dataset of acoustic features or other forms of acoustic features extracted from speech samples of individuals with dysarthria and corresponding severity ratings provided by SLPs. The models learn patterns and relationships between these features and the severity ratings, enabling them to predict/classify the severity of dysarthria in new, unseen cases.

Deep learning, a subset of machine learning, has gained attention in dysarthria severity assessment. Deep Neural Networks (DNN), specifically Long-Short Term Memory (LSTM)[32] and convolutional neural networks (CNNs) [33], have been employed to analyze speech signals and capture intricate temporal and spectral patterns. These models can process raw speech data directly or extract features automatically through multiple layers of computation. They can learn complex representations and make predictions based on the learned patterns, allowing for more accurate severity assessment.

Research studies have demonstrated the potential of these machine learning and deep learning models'

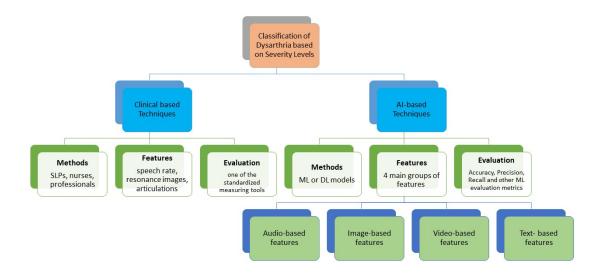


Figure 2: Taxonomy of the study.

potential in objectively quantifying dysarthria's severity. These models offer advantages such as increased efficiency, consistency, and potential for remote or self-administered assessments. However, it is essential to note that these models are still evolving. Their performance may vary depending on factors such as the quality and diversity of training data, feature selection, and the complexity of the dysarthria cases [34][35].

#### 4.2. Features

Machine learning and deep learning models for assessing the severity of dysarthria utilize a range of features extracted from speech signals. These features capture different aspects of speech production and can provide valuable information for severity assessment. We can categorize these features into four main groups: audio-based, image-based, video-based, and text-based. Let's break it down:

## 4.2.1. Audio-based features

Audio-based features for dysarthria classification based on severity can be broadly categorized into acoustic, prosodic, and spectral features [36, 37, 38, 39].

Acoustic features include the fundamental frequency (F0), which corresponds to the pitch of the speech, as well as measures of variability such as jitter and shimmer, reflecting the stability of vocal fold vibration. The Harmonics-to-Noise Ratio (HNR) is another acoustic feature that indicates the ratio of energy in the harmonics of the speech signal to the energy in the noise. Lower HNR values may indicate poor voice quality.

*Prosodic features* play a role in dysarthria classification. Speech rate, which refers to the number of syllables per unit of time, can be affected by the severity of dysarthria. Pause duration, or the length of pauses between speech segments, can provide insights into difficulties in speech production or planning. Additionally,

changes in stress patterns, including variations in intensity, duration, and pitch of syllables, can indicate the presence of dysarthria.

Spectral features offer valuable information for dysarthria classification. Mel-Frequency Cepstral Coefficients (MFCCs) capture the speech signal's short-term power spectrum, providing insights into the speaker's vocal tract shape and function. Linear Predictive Coding (LPC) coefficients represent the spectral envelope of the speech signal, offering information about the vocal tract shape and articulatory movements. Formant frequencies, such as F1, F2, F3, etc., represent the resonant frequencies of the vocal tract during speech production, revealing details about the shape and position of the articulators (tongue, lips, etc.).

The audio-based features mentioned above are not the only types in each group. They represent typical features for dysarthria classification, but additional features can be within each category [26][40]. Table 1 summarizes some studies that used this type of feature for classifying dysarthria based on severity levels or estimating intelligibility for severity classification of dysarthria.

# 4.2.2. Image-based features

Image-based features in the context of dysarthria analysis involve spectrograms, which are visual representations of the frequency content of an audio signal over time.

Spectrograms provide a 2D visual representation of the frequency content of an audio signal over time. They are commonly used in speech and audio analysis, including dysarthria classification. By analyzing spectrograms, various features can be extracted to characterize dysarthric speech. These features include spectral envelope, spectral patterns, spectral variations, and spectral entropy. Spectrograms offer valuable insights into the frequency components and dynamics of speech signals [41].

Mel spectrograms, also known as Mel-frequency spectrograms or Mel-scaled spectrograms, are a type of spectrogram that uses the Mel scale to warp the frequency axis perceptually. The mel scale better aligns with the human auditory system's frequency perception. To generate a mel spectrogram, the audio signal is divided into frames, and the power spectrum is calculated using techniques like the Fast Fourier Transform (FFT). The resulting power spectrum is then transformed to the mel scale using triangular filter banks. Mel spectrograms are particularly useful in speech analysis tasks, including dysarthria classification. They capture essential spectral information, such as formant frequencies and spectral patterns, while aligning with human perception. Mel spectrograms can be used as image-based features for dysarthria analysis or further processed to extract specific characteristics [42]. Additionally, log Mel spectrograms, obtained by taking the logarithm of the values in the Mel spectrogram, can compress the dynamic range and enhance the visualization and analysis of the melfrequency content of the audio signal. [43].

Table 2 summarizes some of the studies in the literature within the collected papers for this review that have utilized image-based features to classify dysarthria based on severity levels.

#### 4.2.3. Video-based features

Video-based features in dysarthria analysis involve analyzing the movement of the lips during speech. Lip Movement Analysis aims to extract visual features that provide information about lip shape, lip dynamics, and lip synchronization. By analyzing the visual cues from lip movements, researchers can gain insights into the articulatory aspects of speech production [44]. However, it is essential to note that classifying the severity of dysarthria solely based on video-based features can be challenging. Dysarthria affects speech mechanisms, including articulation, phonation, resonance, and respiration, which may not be directly observable in video footage. While lip movements can provide valuable information [45], they may not capture the complete picture of dysarthria symptoms.

It is worth mentioning that research and advancements in computer vision techniques, such as facial landmark tracking and optical flow analysis, continue to improve the estimation of video-based features for dysarthria. These techniques enable more precise extraction and interpretation of lip movements and other facial cues, enhancing the potential for video-based research in dysarthria assessment.

### 4.2.4. Text-related features

Text-related features play a crucial role in classifying dysarthria, providing valuable insights into the phonetic characteristics and speech production patterns of individuals with this condition. These features are derived from the analysis of voice signals and focus on the phonetic content of the speech rather than the textual information. Phoneme-level intelligibility serves as a prominent text-related feature, assessing the accuracy of phoneme production by transcribing and comparing

the phonetic content of the speech. Additionally, phonetic distance and articulation errors contribute to the classification process, quantifying the dissimilarity between intended and produced phonemes and identifying specific articulation difficulties. By incorporating these text-related features, dysarthria classification models can capture and analyze the phonetic intricacies of speech, enabling improved understanding and identification of different dysarthria subtypes and severity levels. Such as in [46],[47], and lexical frequency, phonological neighborhood, word class, and lexical familiarity in [48].

#### 4.3. Evaluation

The evaluation techniques used in machine learning and deep learning models play a crucial role in assessing their performance and effectiveness. These techniques provide insights into the model's predictive capabilities and help practitioners make informed decisions[49] [50] [51] [52]. These evaluation metrics include Accuracy measures the overall correctness of the model's predictions. Precision and Recall are commonly used in binary classification tasks, where precision measures the proportion of correctly predicted positive instances and recall measures the proportion of correctly predicted positive instances among all actual positives. The F1Score combines precision and recall into a harmonic mean. Mean Squared Error (MSE) is used for regression tasks and measures the average squared difference between predicted and actual values. Area Under the Curve (AUC) evaluates binary classifiers by calculating the area under the ROC curve. Cross-Validation helps assess model performance across multiple iterations and reduces the risk of overfitting. The Confusion Matrix provides a detailed evaluation of the model's performance by showing true positives, true negatives, false positives, and false negatives. These evaluation techniques assist practitioners in assessing model performance, identifying areas for improvement, and comparing different models for a given task.

## 5. Discussion

Let's analyze these categories of features to compare the effectiveness of different features for classifying dysarthria based on severity levels.

Table 1 summarizes studies that have used audiobased features for dysarthria classification based on severity levels. The studies utilize various audio-based features such as fundamental frequency (F0), jitter, shimmer, harmonics-to-noise ratio (HNR), speech rate, pause duration, stress patterns, Mel-Frequency Cepstral Coefficients (MFCCs), Linear Predictive Coding (LPC) coefficients, and formant frequencies (F1, F2, F3, etc.). Different classification techniques have been used, including Random Forest (RF), Support Vector Machine(SVM), Artificial Neural Network (ANN), Classification and Regression Tree (CART), Naive Bayes (NB), DNN, CNN, LSTM, Residual Networks (ResNet), multi-layer perception(MLP), Extreme Gradient Boosting (XGBoost), Gaussian Mixture Model(GMM), Probabilistic linear discriminant analysis (PLDA), Hidden Markov Model (HMM), and k-nearest neighbor(KNN).

It can be seen that most of the researchers' works were evaluated based on standard publicly available datasets such as TORGO [53], UA Speech [54], Qolt [55], and Numours datasets[56], as well as some other locally collected datasets or other less common foreign languages datasets.

Figure 3 shows the distribution of the included papers chosen for this study based on Categorized Features and Dataset Groups. The table reveals interesting patterns in the distribution of papers across different categorized features and dataset groups in dysarthria classification. Among the featured categories, audio-based techniques demonstrate the highest representation across all dataset groups, including TORGO, UA Speech, Nemours, Qolt, and others. This indicates the prevalent use of audio features in dysarthria classification research across diverse datasets. The prominence of audio-based techniques suggests that researchers prioritize capturing the acoustic characteristics of dysarthric speech for accurate classification.

In contrast, image-based and text-based features have a limited presence, with only a few studies exploring their potential within the UA Speech and other dataset groups. This highlights a potential research gap and suggests the need for further investigation into the utilization of visual information for improved dysarthria classification. Video-based features, on the other hand, are not extensively explored in the selected papers, indicating a less prominent role in dysarthria classification across all dataset groups. We added the mixed category as well, which incorporates multiple feature types, and shows a moderate presence in some dataset groups, underscoring the potential benefits of integrating different modalities to enhance classification performance

The majority of using the audio-based features is due to several advantages. Firstly, they provide a direct measurement of speech by capturing important properties of speech signals, including pitch, intensity, and spectral characteristics. These features are essential in evaluating dysarthria, primarily affecting a person's ability to produce clear and intelligible speech.

Secondly, acoustic features allow for objective and quantitative analysis through signal-processing techniques. This objectivity and quantifiability enhance the consistency and reliability of dysarthria severity assessments. By relying on concrete acoustic properties of speech, these measures provide a more robust evaluation of the condition.

Another benefit is that acoustic analysis is non-invasive and accessible. It can be conducted using standard methods such as recording speech samples with a microphone. This practicality makes acoustic analysis suitable for various settings, including clinics, research studies, and telemedicine. It eliminates the need for invasive procedures or specialized equipment, increasing the ease of implementation.

Furthermore, acoustic features enable longitudinal monitoring of dysarthria progression and treatment outcomes. By analyzing changes in the acoustic properties of speech over time, clinicians and researchers can assess the effectiveness of interventions and track the impact of dysarthria on an individual's communication

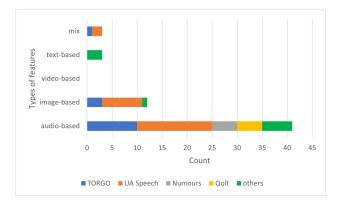


Figure 3: The percentage of each Feature's category

abilities. This longitudinal perspective provides valuable insights into the management and prognosis of dysarthria [57][58][59].

While using acoustic features in dysarthria assessment offers several advantages, there are some limitations to consider. Acoustic features partially represent dysarthric speech, as dysarthria encompasses various communication aspects beyond acoustic properties, such as articulation and prosody. Additionally, acoustic analysis lacks the broader context of communication, failing to capture factors like facial expressions or situational cues that influence speech intelligibility.

The generalizability of acoustic features is limited, as they are derived from controlled environments and may not account for real-world variations in acoustic conditions. The acoustic analysis also primarily focuses on objective measures, potentially overlooking subjective experiences and perceptions of individuals with dysarthria. Moreover, distinguishing between dysarthria subtypes solely based on acoustic analysis can be challenging due to shared acoustic characteristics [60][61]. Despite the promising results of the studies in Table 1, it shows some limitations due to the above-mentioned drawbacks of the audio-based features and other reasons related to the type and size of the dataset used for their training and evaluation and the type of the classifiers used

To overcome these limitations, researchers and clinicians often explore multimodal approaches incorporating other data types, such as images, videos, linguistic features, or perceptual assessments. Integrating multiple modalities can provide a more holistic understanding of dysarthria and improve the accuracy and robustness of the assessment process. As can be seen, some of these attempts are shown in Table 2 for studies that have used image-based features, where most achieved high accuracies around 90% and above. Despite the high performance in most studies in Table 2, still have some limitations, including data scarcity, challenges in collecting dysarthric speech data, lower accuracy due to limited datasets, and confinement to isolated utterances.

Some researchers have derived benefits from employing text-based features within the feature types discussed earlier. These features, extracted from the patient's voice signal and converted into text, have demonstrated promising outcomes. However, their usage

has significant limitations, rendering exclusive reliance on these types impractical. These limitations are represented by the sound quality, which plays a crucial role in the transcription process.

Despite sophisticated speech models employed by Speech-to-Text applications, accuracy can be compromised by the type and specifications of the microphones used. Audio input quality might be degraded if users speak too close to or too far from the microphone, which can subsequently impact the transcription's precision. Another fundamental limitation is the presence of background and environmental noise in the audio input. Nonspeech sounds can interfere with the audio, leading to less accurate transcriptions. The complexity of transcription can also escalate when multiple speakers speak simultaneously or when there is background speech while the primary user is speaking. Specialized vocabulary presents another challenge. Even though Speech-to-Text models can recognize a wide variety of words, they may stumble upon unique terms or industry-specific jargon not included in the model's vocabulary, resulting in potential transcription errors. Accents and dialects within the same language can pose another difficulty. If a speaker's accent deviates substantially from the norm the model has been trained on; transcription accuracy might decline. Finally, language mismatch can considerably affect transcription accuracy. If the language the user speaks differs from what the Speech-to-Text application expects, the transcription will likely be less accurate. For example, if the system is set to transcribe English, but the user speaks Arabic, the output will likely be flawed [62].

Some approaches may combine audio with image or video features to capture a more comprehensive representation of dysarthric speech, such as in [6] where audio-visual joint features are used as input to the CNN and evaluated on UA Speech dataset with an accuracy up to 99.5%. Despite the high performance of their work, their techniques face a few limitations related to manual data pre-processing and high computational power requirements. In [89], the authors applied several types of features rather than mixing them as one input to the classifier to check which type is more effective in classifying dysarthria based on severity levels, where they used Generalized Morse Wavelet (GMW)-based scalogram features (for low-frequency areas) and Mel spectrogram-based features with CNN. They evaluated their work on the UA Speech dataset and achieved the highest accuracy using Scalgram features at 95.17%. Another work where the same has been done in [90] where several experiments on MFCC, the audio combination of features including (ZCR (Zero Crossing Rate), Spectral centroid, spectral roll-off) and Mel spectrogram using KNN, and SVM classifiers and evaluated their work on Torgo with an accuracy of 95% using SVM for mel spectrogram.

It's important to note that the choice of features depends on the specific goals and requirements of the dysarthria severity assessment task. Different features may carry different levels of information and may be more or less suitable for various applications. Researchers and practitioners often select and combine features based

on their relevance and effectiveness in capturing the characteristics of dysarthric speech.

## 5.1. Discussion of Research Question 1

Several types of features were used across the studies: The most common type of feature is the Acoustic feature. These include MFCCs (Mel Frequency Cepstral Coefficients), voice quality, and spectral and prosodic features (such as pitch and rhythm). In a general analysis, most studies used a combination of these features. Some studies focused on features Domain-specific features of speech disorders, such as breathiness features[5] and parameters related to glottal function[63]. Hand-crafted and raw waveform features were used by one study that combined these two approaches[64]. As shown in Figure 4, the mean accuracy obtained by the most utilized classification techniques representing the highest accuracy for each group of the datasets based on audio features reached up to 90s. In specific analysis, Complex, multidimensional feature sets (a combination of features) tend to perform better and are associated with high accuracy rates. For example and referring to individual performance values extracted from the original papers, Ref. [63] using MFCCs, constant-Q cepstral coefficients, prosody, articulation, phonation, and glottal functioning achieved 93.97% accuracy. Another example is the study that achieved 99.9% with UA Speech dataset [65] used a mix of spectral, cepstral, and frame-level features. While combining multiple features and techniques can result in high accuracies, it can also introduce limitations such as increased computational time [63] and difficulties in interpreting feature importance [66].

Prosodic features, including articulation, often demonstrate high performance in dysarthria classification. These features capture essential aspects of speech production, such as rhythm, intonation, and speech clarity, which are crucial for assessing dysarthria severity. In Al-Qatab and Mustafa [68] reported that prosodic features are the best for classifying the mild cases (which are better classified compared to severe and moderate levels due to having more common features among speakers), and the moderate cases over the spectral features, which achieved remarkable results for classifying severe cases. Cepstral features are less effective for classifying severity levels. This is also approved by [67] and Joshy et al. [86].

Some other techniques explored individual features for the classification purposes of ineligibility levels or dysarthria severity levels and discovered some significant features which can achieve a high accuracy rate, such as in [5], which relied on several types of breathiness features represented by Jitter, Shimmer, Harmonicto-Noise Ratio, Harmonic Energy, Harmonic Energy of Residue, Harmonic-to-Signal Ratio, and Glottal-to-Noise Excitation Ratio and reached 96%. In their work, they discovered that Harmonic-related features could distinguish intelligibility levels and achieve the highest accuracy compared to other features but struggled with mild dysarthria. Also, in [77], where the authors rely only on the MFCC features set and achieved up to 99% accuracy rate after experiencing several trails with different utterances lengths, 50, 200, and 300 and found that increasing the length of utterances will lead to better

Table 1
Summary of the related works used the Audio-based features for classifying dysarthria based on the severity levels

Method	Feature	Classifier	Dataset
Chandrashekar	Breathiness features like HNR	SVM	UA Speech and Kan-
et al. [5]			nada Database
Hernandez et al. [31]	prosody, voice quality, MFCC and their combination	RF, SVM, MLP	TORGO, QoLT
Tartarisco et al. [40]	features extracted by VGGish net	ML, TF	Healthy, ataxia subjects
Joshy and Rajan	Acoustic characteristics, articulatory movements, glottal	DNN, CNN, gated re-	UA-Speech and
[63]	functioning, and low-dimensional representations	current units (GRU),and LSTM	TORGO
Yeo et al. [64]	Hand-crafted acoustic features, as well as raw waveforms	Multi task learning (MLT), SVM, MLP, and XGBoost	QoLT
Karjigi et al. [65]	Spectral domain representation features, Cepstral domain representation features, and Frame-level features	PLDA, and ANN	UA speech, TORGO
Yeo et al. [66]	MFCCs, Voice Quality Features, Prosody Features	SVM, ANN	QoLT
Kadi et al. [67]	eleven prosodic features selected by LDA	GMM, SVM	Nemours database
Al-Qatab and Mustafa [68]	four acoustic features: prosody, spectral, cepstral, and voice quality	SVM, LDA, ANN, CART, NB, RF	Nemours database
Bhat et al. [69]	Multi-tapered spectral estimation-based features, Acoustic descriptors for timbre	ANN	UA Speech and TORGO
Hernandez et al. [70]	Standard Prosodic Features and Rhythm-Based Features	RF, SVM, MLP	QoLT, and TORGO
Kachhi et al. [71]	Cepstral Coefficients (TECC, LFCC), and spectral features.	CNN, Light CNN, ResNet	UA Speech
Paja and Falk [72]	Temporal Features, Spectral Features, Voice Quality Features	Mahalanobis distance	UA Speech
Yeo et al. [73]	Language-independent features from diverse speech di- mensions, and language-unique features specific to each language	XGBoost	TORGO, QoLT, and SS- NCE
Kadi et al. [74]	Prosodic features	GMM and SVM	Nemours database
Vyas et al. [75]	prosodic features: (MFCCs), skewness, and formants.	SVM	UA Speech
Kim et al. [76]	10 speech features related to phonetic quality, prosodic quality, and voice quality	SVM	QoLT
Purohit et al. [77]	MFCCs	CNN	UA Speech and home service Corpus
Kachhi et al. [78]	Acoustic/Signal-based Features (SECC, TECC)	CNN, Light CNN, ResNet	UA speech
Mendoza Ramos et al. [79]	10 acoustic features	Discriminant analysis	Dutch speakers and English speakers
Gurugubelli and Vuppala [80]	Single Frequency Filter Bank-based Instantaneous Frequency Cepstral Coefficients (SFFB-IFCC)	I-Vector with PLDA Classifier	UA Speech
Joshy and Rajan [81]	MFCCs and their first two derivatives (a total of 39 features).	SVM, DNN, CNN, LSTM	UA Speech, and TORGO
Narendra and Alku [82]	Glottal Parameters and basic acoustic features	SVM	UA Speech
Gillespie et al. [83]	Spectral, prosodic, Teager Energy Operator (TEO), and glottal waveform features	SVM	UA Speech
Kadi et al. [84]	traditional acoustic features like MFCC and auditory-based cues.	GMM, SVM, and hybrid GMM/SVM	the Nemours database and Torgo
Dahmani et al. [85]	Rhythm Metrics Various rhythm metrics, such as SRVC (speaking rate of VC segments)	SVM	Nemours corpus
Joshy et al. [86]	prosody, articulation, phonation, and glottal	SVM, NB, kNN, and RF	UA Speech and TORGO
Javanmardi et al. [87]	wav2vec features	SVM	UA Speech
DeMino et al. [88]	Frame-level features with global statistics	Forward selection method (FSM)and boosting algorithm	Dataset of 39 dysarthria patients

performance. A similar attempt relying on MFCC is in [91].

Concerning image-based features, the Mel spectrograms and their log or derivatives are the primary types of features and are primarily evaluated on TORGO and UA Speech datasets, as shown in Figure 5. This is a benefit of using only one type to save the preparation and extraction of too many features, as the network is responsible for achieving this task.

The mean accuracies for classifying dysarthria based on severity levels varied across different datasets. Among

the selected datasets, the TORGO dataset exhibited moderate mean accuracies ranging from 73.99% (KNN) to 92.47% (DL), with a notable performance from CNN (88.15%). The UA Speech dataset showed higher mean accuracies, with MLP achieving 98.17% accuracy, followed closely by LSTM (96.94%) and CNN (96.23%). The Qolt dataset had lower mean accuracies, with SVM and RF achieving approximately 70% accuracy. Notably, the Numours dataset demonstrated high mean accuracies of 94.45% (SVM) and 95.8% (RF). The Others dataset had an outstanding performance, with SVM and CNN

Table 2
Summary of the related works used Image-based features for classifying dysarthria based on the severity levels

Method	Feature	Classifier	Dataset
Joshy and Ra- jan [43]	Mel spectrograms	Multi-head attention mechanism (MHA), Multi- task learning, and ResNet	UA Speech
Joshy and Ra- jan [92]	Mel spectrograms	Squeeze-and-excitation (SE) networks	UA Speech
Chandrashekar et al. [93]	perceptually enhanced Fourier trans- form spectrograms and constant-Q transform (CQT) spectrograms	CNN	UA speech and TORGO
Suhas et al. [94]	log Mel spectrograms	CNN	60 patients from each amyotrophic lateral sclerosis (ALS), Parkinson's disease (PD), and healthy controls
Fernández- Díaz and Gallardo- Antolín [95]	Log-mel spectrograms	LSTM networks with Attention mechanism	UA Speech
Montalbo [96]	2D image spectrograms	DySARNet, a densely squeezed-and-excited attention-gated residual neural network	TORGO and UA- Speech
Gupta et al. [97]	Spectrograms of short speech segments	ResNet (CNNs, GMM, and Light CNNs for comparison)	UA Speech
Chandrashekar et al. [98]	Spectrogram and its derivatives	CNN	UA Speech and TORGO

achieving mean accuracies of 96% and 91.8%, respectively.

These mean accuracies provide insights into the effectiveness of classifiers for different datasets in the context of severity-based dysarthria classification. It is important to note that the performance of classifiers can vary depending on the dataset characteristics, such as the number of samples, data quality, and diversity of dysarthria types.

# 5.2. Discussion of Research Question 2

Regarding the classification techniques, advanced algorithms such as Support Vector Machines (SVM), Discriminant analysis, and deep learning models (DNN, CNN, LSTM) were frequently seen and generally delivered a strong performance. For example, SVM was widely utilized across various studies, including [31], [68], [70], [66], [73], [64], [76], [67], [82], [84], [86], and [40].

ANN-based models, seen in studies such as [31], [68], [69], [70], [66], and [65]. RF showed in [31], [68], [70], and [86], and Deep learning techniques (CNN, DNN, LSTM, etc.), applied in studies [63], [71], [81], [87], yielded high accuracies, typically above 90%.

Newer techniques like self-supervised learning with Wav2vec 2.0 XLS-R, as seen in [64], didn't perform as well, achieving only a 65.52% accuracy rate. This highlights the potential challenges in adapting these models for dysarthria severity classification.

From Figure 4, audio-based classification, SVM and RF demonstrated relatively high mean accuracies across most datasets. In the TORGO dataset, SVM achieved a mean accuracy of 77.28%, while RF achieved 79.17% of

mean accuracy. In the UA Speech dataset, both SVM and RF achieved accuracies above 80%, with SVM reaching 80.25% and RF reaching 82.69%. However, in the Qolt dataset, the mean accuracies for SVM and RF dropped to 70% and 70%, respectively. It's important to note that RF achieved a high accuracy of 95.8% in the Numours dataset, while SVM achieved a remarkable accuracy of 96% in the Others dataset.

For image-based classification in Figure 5, CNN and LSTM were the prominent techniques, as deep learning techniques is a direct way to extract features and classify them into several severity levels using their different layers functions. In the TORGO dataset, CNN achieved an accuracy of 88.15%, while LSTM achieved 99.2% accuracy. In the UA Speech dataset, CNN demonstrated high accuracy of 96.23%, and LSTM achieved 96.94% accuracy. Notably, ResNet only yields accuracy for the UA Speech datasets around 97%.

These results suggest that for audio-based classification, SVM and RF show consistent performance across various datasets. Despite the effectiveness of using SVM, it is reported in Kadi et al. [84]that it is not adequate to process utterances with diverse time lengths and needs to unify the time lengths for all speech utterances. For image-based classification, CNN and LSTM exhibit high accuracies. However, further investigation and evaluation are needed to determine these techniques' generalizability and performance on more extensive and diverse dysarthria datasets.

The mean accuracies presented in the tables were obtained by calculating the average of the highest accuracies reported in the original papers for each dataset

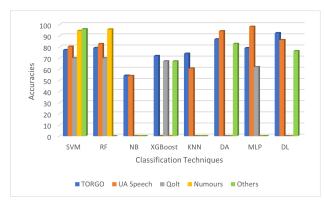


Figure 4: Mean Accuracy for Audio-based Features for common datasets

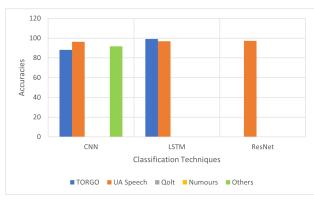


Figure 5: Mean Accuracy for Image-based Features for common datasets

group. To ensure accuracy and reliability, only the highest accuracy value from each paper was considered. This approach allows for a comprehensive analysis of the performance of the classification techniques across different datasets, taking into account the number of papers that reported the respective accuracies. By aggregating the highest accuracies, we provide a representative measure of the overall performance of the techniques for each dataset group, considering the number of papers that contributed to calculating the mean accuracies.

# 6. Limitations

The audio-based features used in dysarthria classification have certain limitations, which can be compared to image-based features. By analyzing these limitations, we can gain insights into the techniques and accuracies associated with each feature type.

Regarding audio-based features, one limitation mentioned in [5] is the struggle to accurately classify cases with mild dysarthria. This implies that audio-based features might be more effective in capturing severe dysarthria characteristics, while their discriminatory power decreases with milder cases. Additionally, there are difficulties in interpreting the features and handling data diversity [31]. This suggests that some audio-based features may lack clear interpretability and may not be robust enough to handle the diverse dysarthria characteristics. Furthermore, there are limitations related to small datasets and the focus on a single type of dysarthria [68], which can hinder the generalizability of the findings.

On the other hand, image-based features also have their limitations. For instance, inter-subject variability and gender bias can impact the performance of image-based models [92]. These limitations suggest that certain image-based features might be more influenced by individual differences, making them less reliable for general dysarthria classification. Another limitation is the computational time required for specific image-based techniques [93], which can hinder their real-time application.

When comparing the techniques and accuracies, it is essential to note that the performance varies across different studies and datasets. Audio-based features have reported accuracies ranging from 40.41% to 95.80% [68], whereas image-based features have achieved accuracies ranging from 72% to 99.20% [93, 96]. These variations highlight the influence of the chosen techniques and datasets on the achieved accuracies.

In summary, both audio-based and image-based features are limited in dysarthria classification. Audio-based features struggle with mild dysarthria cases, interpretation difficulties, and limited datasets, while image-based features face challenges related to inter-subject variability, gender bias, and computational time. These limitations influence the techniques used with these features and the resulting accuracies, emphasizing the need to carefully consider the feature type, technique selection, and dataset characteristics in dysarthria classification research.

## 7. Suggested solutions

To overcome the limitation cited above, some analyses are performed using computer vision features that can improve the performance of each method.

Solving the issue of larger data size: Instead of utilizing only the standard machine learning techniques or the deep learning techniques which require large amounts of the dataset, we suggest using rule-based models such as Adaptive Neuro-Fuzzy Inference System (ANFIS), which outperform the other methods in terms of the necessity to small data size for training either in terms of features or samples, the more interpretable model explicitly compared to deep learning, adaptable and generalized model due to the generated rules which can be adjusted based on the specific application and the used data nature and the non-linearity and noise handling [99][100].

Some new features are suggested to be used by researchers in this field inclduing Nonlinear Dynamical Features, Vocal Tract Resonance Estimation, Source-Filter Separation Features, Microprosody Analysis, and Articulatory Kinematics.

Nonlinear Dynamical Features: Dysarthria affects the coordination and control of speech production, resulting in altered dynamics. Nonlinear dynamical features, such as recurrence plots(can help visualize temporal patterns in speech features that could indicate varying levels of severity [101]), Variations in the fractal dimension of speech signals could distinguish between different severity levels of dysarthria [102], or entropy measures

(e.g., permutation entropy, sample entropy) (lower entropy indicates a less complex, more predictable one, and dysarthria's speech signals may become less complicated due to muscle weakness and coordination difficulties, can capture the underlying complexity and temporal organization of dysarthric speech signals [103]). We know these features have not been utilized to classify dysarthria. Still, it is used with speech analysis, such as in [104], and it is helpful to be calculated and combined with video-based features.

Vocal Tract Resonance Estimation: Investigate techniques for estimating the vocal tract resonances directly from the speech signal. This could involve formant tracking, subspace-based analysis, or model-based estimation. This feature type was utilized for classifying dysarthria subtypes in [105]. These features alone may not be sufficient for classifying dysarthria based on severity levels. They primarily capture information about the resonant characteristics of the vocal tract, which can be influenced by various factors such as vocal fold movement, articulatory precision, and vocal tract shape. So, combining them with other acoustic features may help build a successful classification model.

Source-Filter Separation Features: Dysarthria can affect the speaker's source (glottal excitation) and filter (vocal tract) components. Extracting features that specifically focus on the characteristics of the glottal source, such as measures related to the glottal flow derivative or instantaneous frequency, in combination with standard spectral features, can provide a comprehensive representation for dysarthria classification. These features have been utilized in [106] for automatic dysarthric speech recognition and evaluated on the TORGO dataset, which contains diverse severity samples. Still, there is no specific classification model based on these features.

**Microprosody Analysis**: Investigate features that capture fine-grained temporal variations within short speech segments, called microporosity. This could involve analyzing pitch fluctuations, intensity modulations, or spectral changes at very short time scales (e.g., using wavelet transforms or time-frequency representations). These features can potentially capture subtle dysarthric speech dynamics. These types of features used mainly with speech synthesis and its variations were also explored in many studies such as [107].

Articulatory Kinematics: Consider exploring articulatory and kinematic features, such as lip or jaw motion tracking. Using techniques like optical motion capture or electromagnetic articulography, features can be extracted that describe articulatory gestures' spatial and temporal characteristics, which may be informative for dysarthria classification. These features have a notable correlation with the severity levels of dysarthria as discussed in [108], and they have been used to classify vocal segments into control, PD, and ALS, but still no specific work on classifying dysarthria based on severity utilizing this type of features.

# 8. Conclusion

Accurately classifying dysarthria based on severity levels is crucial for clinicians and researchers, as

it provides valuable insights into the disorder's impact on communication abilities. This classification enables tailored treatment strategies, objective assessment measures, and progress tracking, ultimately improving the quality of life for individuals with dysarthria. This review aimed to analyze the classification of dysarthria based on severity levels, focusing on the effectiveness of different features and artificial intelligence techniques. The use of audio-based features, mainly acoustic analysis, has been the predominant approach in dysarthria classification. However, limitations such as data scarcity and challenges in real-time implementation have been observed. To overcome these limitations, multimodal systems integrating other data types have been explored. Complex, multidimensional feature sets tend to perform better but can introduce challenges such as increased computational time. The choice of AI techniques does not directly correlate with higher accuracy, as both traditional and advanced techniques have shown varying performance. In conclusion, integrating multiple modalities, exploring new feature types, and carefully selecting AI techniques can improve the accuracy and efficiency of dysarthria classification systems. Future research should address the identified limitations and refine the classification methods for enhanced clinical utility.

# **Declaration of competing interest**

The authors declare no conflicts of interest.

# Acknowledgements

This publication was made by International Research Collaboration Co-Fund (IRCC) Cycle 06 (2023-2024) No. IRCC-2023-223 from the Qatar National Research Fund (a member of the Qatar Foundation). The statements made herein are solely the responsibility of the authors.

#### References

- L. D. Shriberg, E. A. Strand, K. J. Jakielski, H. L. Mabie, Estimates of the prevalence of speech and motor speech disorders in persons with complex neurodevelopmental disorders, Clinical Linguistics & Phonetics 33 (2019) 707–736.
- [2] J. R. Duffy, Motor speech disorders e-book: Substrates, differential diagnosis, and management, Elsevier Health Sciences, 2019.
- [3] R. D. Kent, J. C. Rosenbek, Acoustic patterns of apraxia of speech, Journal of Speech, Language, and Hearing Research 26 (1983) 231–249.
- [4] M. Laganaro, C. Fougeron, M. Pernon, N. Levêque, S. Borel, M. Fournet, S. Catalano Chiuvé, U. Lopez, R. Trouville, L. Ménard, et al., Sensitivity and specificity of an acoustic-and perceptual-based tool for assessing motor speech disorders in french: The monpage-screening protocol, Clinical Linguistics & Phonetics 35 (2021) 1060–1075.
- [5] H. Chandrashekar, V. Karjigi, N. Sreedevi, Breathiness indices for classification of dysarthria based on type and speech intelligibility, in: 2019 International Conference on Wireless Communications Signal Processing and Networking (WiSPNET), IEEE, 2019, pp. 266–270.
- [6] H. Tong, Automatic assessment of dysarthric severity level using audio-video cross-modal approach in deep learning, Master's thesis, 2020.
- [7] G. Noffs, T. Perera, S. C. Kolbe, C. J. Shanahan, F. M. Boonstra, A. Evans, H. Butzkueven, A. van der Walt, A. P. Vogel, What

- speech can tell us: A systematic review of dysarthria characteristics in multiple sclerosis, Autoimmunity reviews 17 (2018) 1202–1209.
- [8] P. Gandhi, S. Tobin, M. Vongphakdi, A. Copley, K. Watter, A scoping review of interventions for adults with dysarthria following traumatic brain injury, Brain injury 34 (2020) 466– 479.
- [9] C. Mackenzie, Dysarthria in stroke: a narrative review of its description and the outcome of intervention, International journal of speech-language pathology 13 (2011) 125–136.
- [10] J. Lee, A. Madhavan, E. Krajewski, S. Lingenfelter, Assessment of dysarthria and dysphagia in patients with amyotrophic lateral sclerosis: Review of the current evidence, Muscle & Nerve 64 (2021) 520–531.
- [11] E. Finch, A. F. Rumbach, S. Park, Speech pathology management of non-progressive dysarthria: a systematic review of the literature, Disability and Rehabilitation 42 (2020) 296–306.
- [12] C. Whillans, M. Lawrie, E. A. Cardell, C. Kelly, R. Wenke, A systematic review of group intervention for acquired dysarthria in adults, Disability and Rehabilitation 44 (2022) 3002–3018.
- [13] Y.-J. Park, J.-M. Lee, Effect of acupuncture intervention and manipulation types on poststroke dysarthria: a systematic review and meta-analysis, Evidence-Based Complementary and Alternative Medicine 2020 (2020).
- [14] A. K. Shanmugam, R. Marimuthu, A critical analysis and review of assistive technology: advancements, laws, and impact on improving the rehabilitation of dysarthric patients, Handbook of Decision Support Systems for Neurological Disorders (2021) 263–281.
- [15] H. P. Rowe, S. E. Gutz, M. F. Maffei, K. Tomanek, J. R. Green, Characterizing dysarthria diversity for automatic speech recognition: A tutorial from the clinical perspective, Frontiers in Computer Science (2022) 43.
- [16] R. Chiaramonte, M. Vecchio, A systematic review of measures of dysarthria severity in stroke patients, Pm&r 13 (2021) 314– 324
- [17] R. Palmer, P. Enderby, Methods of speech therapy treatment for stable dysarthria: A review, Advances in Speech Language Pathology 9 (2007) 140–153.
- [18] A. Huang, K. Hall, C. Watson, S. R. Shahamiri, A review of automated intelligibility assessment for dysarthric speakers, in: 2021 International Conference on Speech Technology and Human-Computer Dialogue (SpeD), IEEE, 2021, pp. 19–24.
- [19] K. L. Stipancic, K. M. Palmer, H. P. Rowe, Y. Yunusova, J. D. Berry, J. R. Green, "you say severe, i say mild": Toward an empirical classification of dysarthria severity, Journal of Speech, Language, and Hearing Research 64 (2021) 4718–4735.
- [20] M. R. McNeil, K. J. Ballard, J. R. Duffy, J. Wambaugh, P. van Lieshout, B. Maassen, H. Terband, Apraxia of speech theory, assessment, differential diagnosis, and treatment: Past, present, and future, Speech motor control in normal and disordered speech: Future developments in theory and methodology (2017) 195–221.
- [21] J. M. Barkmeier-Kraemer, H. M. Clark, Speech-language pathology evaluation and management of hyperkinetic disorders affecting speech and swallowing function, Tremor and Other Hyperkinetic Movements 7 (2017).
- [22] P. Enderby, Frenchay dysarthria assessment, British Journal of Disorders of Communication 15 (1980) 165–173.
- [23] K. M. Yorkston, D. R. Beukelman, C. Traynor, Assessment of intelligibility of dysarthric speech, Pro-ed Austin, TX, 1984.
- [24] S. J. Robertson, Dysarthria profile, Communication Skill Builders, 1987.
- [25] B. H. Jacobson, A. Johnson, C. Grywalski, A. Silbergleit, G. Jacobson, M. S. Benninger, C. W. Newman, The voice handicap index (vhi) development and validation, American journal of speech-language pathology 6 (1997) 66–70.
- [26] D. B. Freed, Motor speech disorders: diagnosis and treatment, Plural Publishing, 2018.
- [27] S. Blaustein, A. Bar, Reliability of perceptual voice assessment, Journal of communication disorders 16 (1983) 157–161.
- [28] S. A. S. Al Yaari, N. Almaflehi, Validation of communication activities of daily living-(cadl-2) on arab aphasics: Controlled study, J Study Engl Linguist 2 (2014) 34.

- [29] A. M. Altaher, S. Y. Chu, R. A. Razak, et al., A report of assessment tools for individuals with dysarthria, The Open Public Health Journal 12 (2019).
- [30] J. I. Godino-Llorente, P. Gomez-Vilda, Automatic detection of voice impairments by means of short-term cepstral parameters and neural network based detectors, IEEE Transactions on Biomedical Engineering 51 (2004) 380–384.
- [31] A. Hernandez, S. Kim, M. Chung, Prosody-based measures for automatic severity assessment of dysarthric speech, Applied Sciences 10 (2020) 6999.
- [32] A. Gallardo-Antolín, J. M. Montero, An auditory saliency pooling-based lstm model for speech intelligibility classification, Symmetry 13 (2021) 1728.
- [33] J. C. Vásquez-Correa, T. Arias-Vergara, J. R. Orozco-Arroyave, E. Nöth, A multitask learning approach to assess the dysarthria severity in patients with parkinson's disease., in: INTER-SPEECH, 2018, pp. 456–460.
- [34] L. Parisi, N. RaviChandran, M. L. Manaog, Feature-driven machine learning to improve early diagnosis of parkinson's disease, Expert Systems with Applications 110 (2018) 182–190.
- [35] L. Xu, J. Liss, V. Berisha, Dysarthria detection based on a deep learning model with a clinically-interpretable layer, JASA Express Letters 3 (2023) 015201.
- [36] R. D. Kent, G. Weismer, J. F. Kent, H. K. Vorperian, J. R. Duffy, Acoustic studies of dysarthric speech: Methods, progress, and potential, Journal of communication disorders 32 (1999) 141– 186
- [37] S. Sapir, L. O. Ramig, J. L. Spielman, C. Fox, Formant centralization ratio: A proposal for a new acoustic measure of dysarthric speech (2010).
- [38] H. Kim, M. Hasegawa-Johnson, A. Perlman, Vowel contrast and speech intelligibility in dysarthria, Folia Phoniatrica et Logopaedica 63 (2011) 187–194.
- [39] J. R. Orozco-Arroyave, J. D. Arias-Londoño, J. F. Vargas-Bonilla, M. C. Gonzalez-Rátiva, E. Nöth, New spanish speech corpus database for the analysis of people suffering from parkinson's disease., in: LREC, 2014, pp. 342–347.
- [40] G. Tartarisco, R. Bruschetta, S. Summa, L. Ruta, M. Favetta, M. Busa, A. Romano, E. Castelli, F. Marino, A. Cerasa, et al., Artificial intelligence for dysarthria assessment in children with ataxia: A hierarchical approach, IEEE Access 9 (2021) 166720– 166735.
- [41] D. Kempler, D. Van Lancker, Effect of speech task on intelligibility in dysarthria: A case study of parkinson's disease, Brain and language 80 (2002) 449–464.
- [42] U. Ayvaz, H. Gürüler, F. Khan, N. Ahmed, T. Whangbo, A. Bobomirzaevich, Automatic speaker recognition using mel-frequency cepstral coefficients through machine learning, CMC-Computers Materials & Continua 71 (2022).
- [43] A. A. Joshy, R. Rajan, Dysarthria severity classification using multi-head attention and multi-task learning, Speech Communication 147 (2023) 1–11.
- [44] A. Bandini, J. R. Green, B. Taati, S. Orlandi, L. Zinman, Y. Yunusova, Automatic detection of amyotrophic lateral sclerosis (als) from video-based analysis of facial movements: speech and non-speech tasks, in: 2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018), IEEE, 2018, pp. 150–157.
- [45] H. Ackermann, I. Hertrich, G. Scharf, Kinematic analysis of lower lip movements in ataxic dysarthria, Journal of Speech, Language, and Hearing Research 38 (1995) 1252–1259.
- [46] W. Xue, R. van Hout, C. Cucchiarini, H. Strik, Assessing speech intelligibility of pathological speech in sentences and word lists: The contribution of phoneme-level measures, Journal of Communication Disorders 102 (2023) 106301.
- [47] W. Xue, R. van Hout, C. Cucchiarini, H. Strik, Assessing speech intelligibility of pathological speech: test types, ratings and transcription measures, Clinical Linguistics & Phonetics 37 (2023) 52–76.
- [48] K. Lehner, W. Ziegler, The impact of lexical and articulatory factors in the automatic selection of test materials for a webbased assessment of intelligibility in dysarthria, Journal of Speech, Language, and Hearing Research 64 (2021) 2196–2212.

- [49] N. Japkowicz, M. Shah, Evaluating learning algorithms: a classification perspective, Cambridge University Press, 2011.
- [50] D. M. Powers, Evaluation: from precision, recall and f-measure to roc, informedness, markedness and correlation, arXiv preprint arXiv:2010.16061 (2020).
- [51] A. P. Bradley, The use of the area under the roc curve in the evaluation of machine learning algorithms, Pattern recognition 30 (1997) 1145–1159.
- [52] T. Hastie, R. Tibshirani, J. H. Friedman, J. H. Friedman, The elements of statistical learning: data mining, inference, and prediction, volume 2, Springer, 2009.
- [53] F. Rudzicz, A. K. Namasivayam, T. Wolff, The torgo database of acoustic and articulatory speech from speakers with dysarthria, Language Resources and Evaluation 46 (2012) 523–541.
- [54] H. Kim, M. Hasegawa-Johnson, A. Perlman, J. Gunderson, T. S. Huang, K. Watkin, S. Frame, Dysarthric speech database for universal access research, in: Ninth Annual Conference of the International Speech Communication Association, 2008.
- [55] D.-L. Choi, B.-W. Kim, Y.-W. Kim, Y.-J. Lee, Y. Um, M. Chung, Dysarthric speech database for development of qolt software technology., in: LREC, 2012, pp. 3378–3381.
- [56] X. Menendez-Pidal, J. B. Polikoff, S. M. Peters, J. E. Leonzio, H. T. Bunnell, The nemours database of dysarthric speech, in: Proceeding of Fourth International Conference on Spoken Language Processing. ICSLP'96, volume 3, IEEE, 1996, pp. 1962–1965.
- [57] R. Kent, K. Hustad, Speech production: development (2009).
- [58] K. López-de Ipiña, J.-B. Alonso, C. M. Travieso, J. Solé-Casals, H. Egiraun, M. Faundez-Zanuy, A. Ezeiza, N. Barroso, M. Ecay-Torres, P. Martinez-Lage, et al., On the selection of noninvasive methods based on speech analysis oriented to automatic alzheimer disease diagnosis, Sensors 13 (2013) 6730–6745.
- [59] S. Luz, Longitudinal monitoring and detection of alzheimer's type dementia from spontaneous speech data, in: 2017 IEEE 30th International Symposium on Computer-Based Medical Systems (CBMS), IEEE, 2017, pp. 45–46.
- [60] G. S. Turner, K. Tjaden, G. Weismer, The influence of speaking rate on vowel space and speech intelligibility for individuals with amyotrophic lateral sclerosis, Journal of Speech, Language, and Hearing Research 38 (1995) 1001–1013.
- [61] G. Weismer, J.-Y. Jeng, J. S. Laures, R. D. Kent, J. F. Kent, Acoustic and intelligibility characteristics of sentence production in neurogenic speech disorders, Folia Phoniatrica et Logopaedica 53 (2001) 1–18.
- [62] S. Modale, G. Sable, R. Deshmukh, A review: Devnagri speech to text for marathwada region, in: Proceedings of International Conference on Wireless Communication: ICWiCOM 2019, Springer, 2020, pp. 525–532.
- [63] A. A. Joshy, R. Rajan, Automated dysarthria severity classification: A study on acoustic features and deep learning techniques, IEEE Transactions on Neural Systems and Rehabilitation Engineering 30 (2022) 1147–1157.
- [64] E. J. Yeo, K. Choi, S. Kim, M. Chung, Automatic severity classification of dysarthric speech by using self-supervised model with multi-task learning, in: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2023, pp. 1–5.
- [65] V. Karjigi, N. Sreedevi, et al., Speech intelligibility assessment of dysarthria using fisher vector encoding, Computer Speech & Language 77 (2023) 101411.
- [66] E. J. Yeo, S. Kim, M. Chung, Automatic severity classification of korean dysarthric speech using phoneme-level pronunciation features., in: Interspeech, 2021, pp. 4838–4842.
- [67] K. L. Kadi, S. A. Selouani, B. Boudraa, M. Boudraa, Automated diagnosis and assessment of dysarthric speech using relevant prosodic features, in: Transactions on Engineering Technologies: Special Volume of the World Congress on Engineering 2013, Springer, 2014, pp. 529–542.
- [68] B. A. Al-Qatab, M. B. Mustafa, Classification of dysarthric speech according to the severity of impairment: an analysis of acoustic features, IEEE Access 9 (2021) 18183–18194.
- [69] C. Bhat, B. Vachhani, S. K. Kopparapu, Automatic assessment of dysarthria severity level using audio descriptors, in: 2017 IEEE International Conference on Acoustics, Speech and Signal

- Processing (ICASSP), IEEE, 2017, pp. 5070-5074.
- [70] A. Hernandez, E. J. Yeo, S. Kim, M. Chung, Dysarthria detection and severity assessment using rhythm-based metrics., in: INTERSPEECH, 2020, pp. 2897–2901.
- [71] A. Kachhi, A. Therattil, A. T. Patil, H. B. Sailor, H. A. Patil, Teager energy cepstral coefficients for classification of dysarthric speech severity-level, in: 2022 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), IEEE, 2022, pp. 1462–1468.
- [72] M. S. Paja, T. H. Falk, Automated dysarthria severity classification for improved objective intelligibility assessment of spastic dysarthric speech, in: Thirteenth Annual Conference of the International Speech Communication Association, 2012.
- [73] E. J. Yeo, K. Choi, S. Kim, M. Chung, Cross-lingual dysarthria severity classification for english, korean, and tamil, in: 2022 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), IEEE, 2022, pp. 566–574.
- [74] K. Kadi, S. Selouani, B. Boudraa, M. Boudraa, Discriminative prosodic features to assess the dysarthria severity levels, in: Proceedings of the World Congress on Engineering, volume 3, 2013
- [75] G. Vyas, M. K. Dutta, J. Prinosil, P. Harár, An automatic diagnosis and assessment of dysarthric speech using speech disorder specific prosodic features, in: 2016 39th International Conference on Telecommunications and Signal Processing (TSP), IEEE, 2016, pp. 515–518.
- [76] M. J. Kim, J. Yoo, H. Kim, Dysarthric speech recognition using dysarthria-severity-dependent and speaker-adaptive models., in: Interspeech, 2013, pp. 3622–3626.
- [77] M. Purohit, M. Parmar, M. Patel, H. Malaviya, H. A. Patii, Weak speech supervision: A case study of dysarthria severity classification, in: 2020 28th European Signal Processing Conference (EUSIPCO), IEEE, 2021, pp. 101–105.
- [78] A. Kachhi, A. Therattil, A. T. Patil, H. B. Sailor, H. A. Patil, Significance of energy features for severity classification of dysarthria, in: Speech and Computer: 24th International Conference, SPECOM 2022, Gurugram, India, November 14–16, 2022, Proceedings, Springer, 2022, pp. 325–337.
- [79] V. Mendoza Ramos, A. Lowit, L. Van den Steen, H. A. Kairuz Hernandez-Diaz, M. E. Hernandez-Diaz Huici, M. De Bodt, G. Van Nuffelen, Acoustic identification of sentence accent in speakers with dysarthria: cross-population validation and severity related patterns, Brain Sciences 11 (2021) 1344.
- [80] K. Gurugubelli, A. K. Vuppala, Analytic phase features for dysarthric speech detection and intelligibility assessment, Speech Communication 121 (2020) 1–15.
- [81] A. A. Joshy, R. Rajan, Automated dysarthria severity classification using deep learning frameworks, in: 2020 28th European Signal Processing Conference (EUSIPCO), IEEE, 2021, pp. 116–120
- [82] N. P. Narendra, P. Alku, Automatic intelligibility assessment of dysarthric speech using glottal parameters, Speech Communication 123 (2020) 1–9.
- [83] S. Gillespie, Y.-Y. Logan, E. Moore, J. Laures-Gore, S. Russell, R. Patel, Cross-database models for the classification of dysarthria presence., in: Interspeech, 2017, pp. 3127–3131.
- [84] K. L. Kadi, S. A. Selouani, B. Boudraa, M. Boudraa, Fully automated speaker identification and intelligibility assessment in dysarthria disease using auditory knowledge, Biocybernetics and Biomedical Engineering 36 (2016) 233–247.
- [85] H. Dahmani, S.-A. Selouani, N. Doghmane, D. O'Shaughnessy, M. Chetouani, On the relevance of using rhythmic metrics and svm to assess dysarthric severity, International Journal of Biometrics 6 (2014) 248–271.
- [86] A. A. Joshy, P. Parameswaran, S. R. Nair, R. Rajan, Statistical analysis of speech disorder specific features to characterise dysarthria severity level, in: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2023, pp. 1–5.
- [87] F. Javanmardi, S. Tirronen, M. Kodali, S. R. Kadiri, P. Alku, Wav2vec-based detection and severity level classification of dysarthria from speech, in: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing

- (ICASSP), IEEE, 2023, pp. 1-5.
- [88] A. DeMino, R. Kubichek, K. Caves, Assessing dysarthria severity using global statistics and boosting, in: 2011 Conference Record of the Forty Fifth Asilomar Conference on Signals, Systems and Computers (ASILOMAR), IEEE, 2011, pp. 1103–1106.
- [89] A. Kachhi, A. Therattil, P. Gupta, H. A. Patil, Continuous wavelet transform for severity-level classification of dysarthria, in: Speech and Computer: 24th International Conference, SPECOM 2022, Gurugram, India, November 14–16, 2022, Proceedings, Springer, 2022, pp. 312–324.
- [90] S. A. Kumar, T. Keerthivasan, S. Sasikala, et al., Towards improving the performance of dysarthric speech severity assessment system, in: 2022 International Conference on Computer Communication and Informatics (ICCCI), IEEE, 2022, pp. 1–6.
- [91] S. H. Lee, M. Kim, H. G. Seo, B.-M. Oh, G. Lee, J.-H. Leigh, Assessment of dysarthria using one-word speech recognition with hidden markov models, Journal of Korean medical science 34 (2019).
- [92] A. A. Joshy, R. Rajan, Dysarthria severity assessment using squeeze-and-excitation networks, Biomedical Signal Processing and Control 82 (2023) 104606.
- [93] H. Chandrashekar, V. Karjigi, N. Sreedevi, Investigation of different time-frequency representations for intelligibility assessment of dysarthric speech, Ieee transactions on neural systems and rehabilitation engineering 28 (2020) 2880–2889.
- [94] B. Suhas, J. Mallela, A. Illa, B. Yamini, N. Atchayaram, R. Yadav, D. Gope, P. K. Ghosh, Speech task based automatic classification of als and parkinson's disease and their severity using log mel spectrograms, in: 2020 international conference on signal processing and communications (SPCOM), IEEE, 2020, pp. 1–5.
- [95] M. Fernández-Díaz, A. Gallardo-Antolín, An attention long short-term memory based system for automatic classification of speech intelligibility, Engineering Applications of Artificial Intelligence 96 (2020) 103976.
- [96] F. J. Montalbo, Dysarnet: A densely squeezed-and-excited attention-gated residual deep learning model for dysarthric speech recognition and severity estimation, Available at SSRN 4442941 (January 2023).
- [97] S. Gupta, A. T. Patil, M. Purohit, M. Parmar, M. Patel, H. A. Patil, R. C. Guido, Residual neural network precisely quantifies dysarthria severity-level based on short-duration speech segments. Neural Networks 139 (2021) 105–117.
- [98] H. Chandrashekar, V. Karjigi, N. Sreedevi, Spectro-temporal representation of speech for intelligibility assessment of dysarthria, IEEE Journal of Selected Topics in Signal Processing 14 (2019) 390–399.
- [99] S. Akkoç, An empirical comparison of conventional techniques, neural networks and the three stage hybrid adaptive neuro fuzzy inference system (anfis) model for credit scoring analysis: The case of turkish credit card data, European Journal of Operational Research 222 (2012) 168–178.
- [100] A. Subasi, Application of adaptive neuro-fuzzy inference system for epileptic seizure detection using wavelet feature extraction, Computers in biology and medicine 37 (2007) 227– 244
- [101] O. Geman, Data processing for parkinson's disease: Tremor, speech and gait signal analysis, in: 2011 E-Health and Bioengineering Conference (EHB), IEEE, 2011, pp. 1–4.
- [102] A. P. Accardo, E. Mumolo, An algorithm for the automatic differentiation between the speech of normals and patients with friedreich's ataxia based on the short-time fractal dimension, Computers in biology and medicine 28 (1998) 75–89.
- [103] F. Rudzicz, Towards a noisy-channel model of dysarthria in speech recognition, in: Proceedings of the NAACL HLT 2010 Workshop on Speech and Language Processing for Assistive Technologies, 2010, pp. 80–88.
- [104] E. S. Jackson, M. Tiede, M. A. Riley, D. Whalen, Recurrence quantification analysis of sentence-level speech kinematics, Journal of Speech, Language, and Hearing Research 59 (2016) 1315–1326.
- [105] N. Lévêque, A. Slis, L. Lancia, G. Bruneteau, C. Fougeron, Acoustic change over time in spastic and/or flaccid dysarthria

- in motor neuron diseases, Journal of Speech, Language, and Hearing Research 65 (2022) 1767–1783.
- [106] Z. Yue, E. Loweimi, Z. Cvetkovic, Raw source and filter modelling for dysarthric speech recognition, in: ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2022, pp. 7377–7381.
- [107] M. Vainio, T. Altosaar, Modeling the microprosody of pitch and loudness for speech synthesis with neural networks., in: ICSLP, 1998.
- [108] J. Lee, M. Bell, Z. Simmons, Articulatory kinematic characteristics across the dysarthria severity spectrum in individuals with amyotrophic lateral sclerosis, American Journal of Speech-Language Pathology 27 (2018) 258–269.